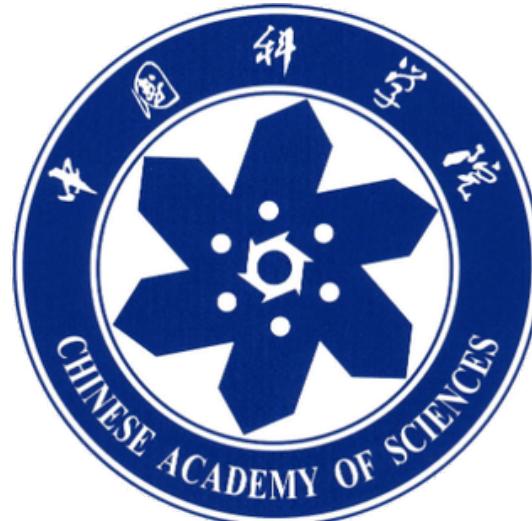
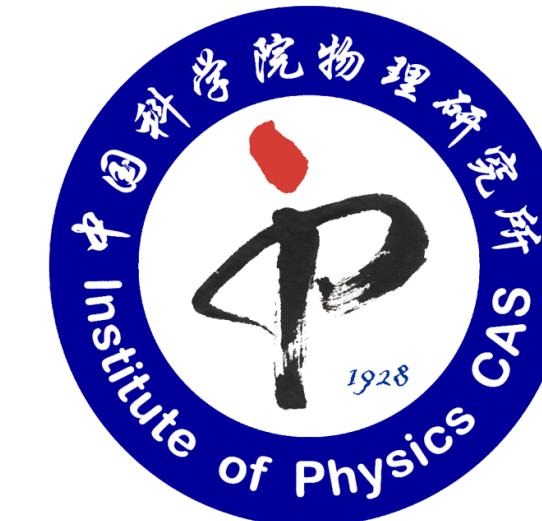


Generative models for physicists

Lei Wang (王磊)

Institute of Physics, CAS

<https://wangleiphy.github.io>



Plan

①

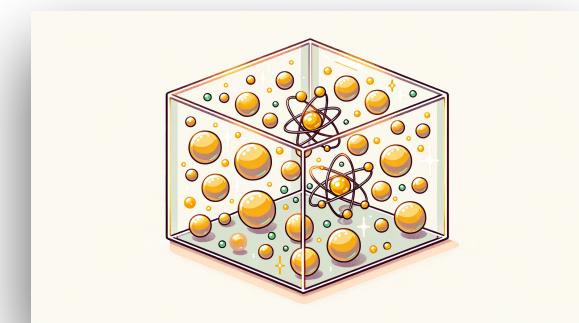
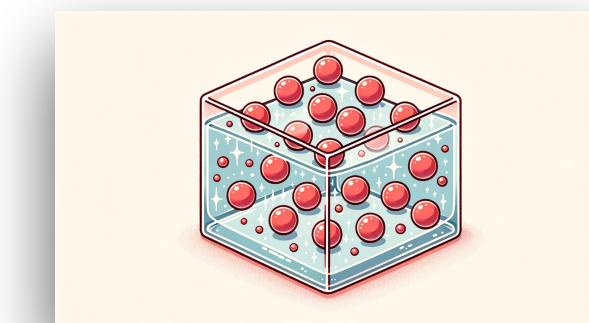
Motivations

②

Generative models and their physics genes

③

Applications: electron gases and dense hydrogen



Science is more than fitting, so is machine learning

Discriminative learning



Generative learning



$$y = f(x)$$

or $p(y | x)$

$$p(x, y)$$

Science is more than fitting, so is machine learning

Discriminative learning



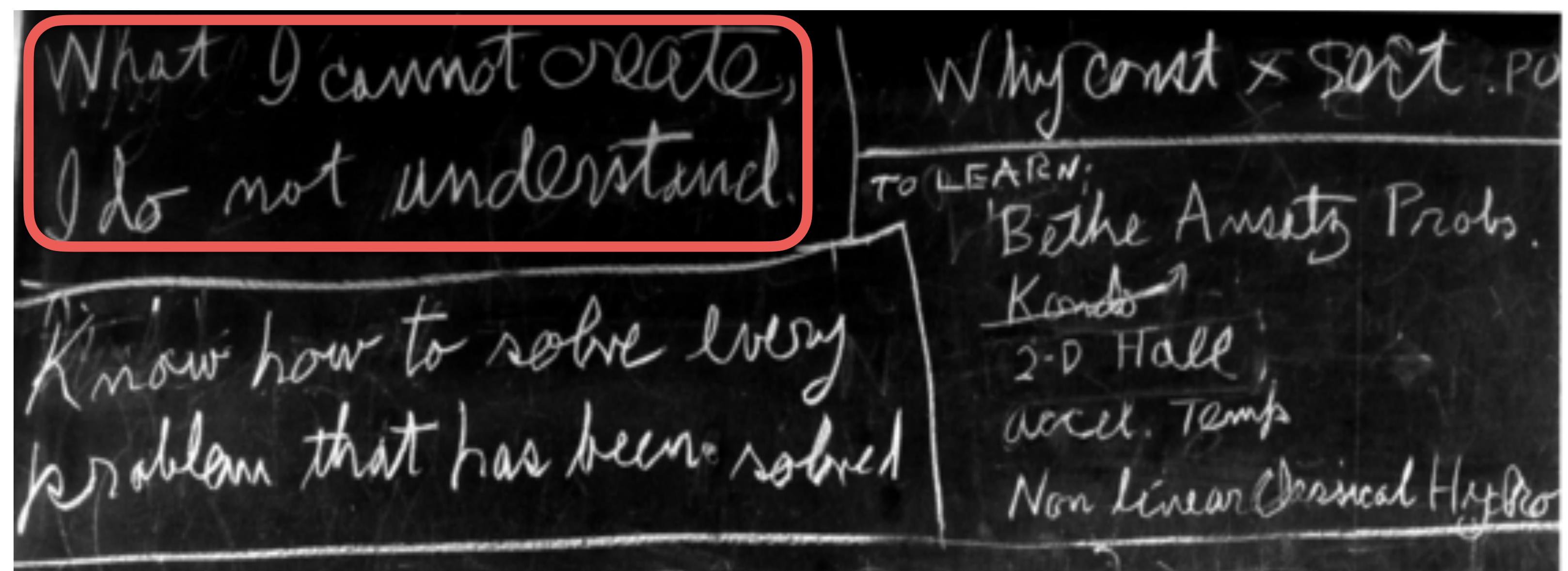
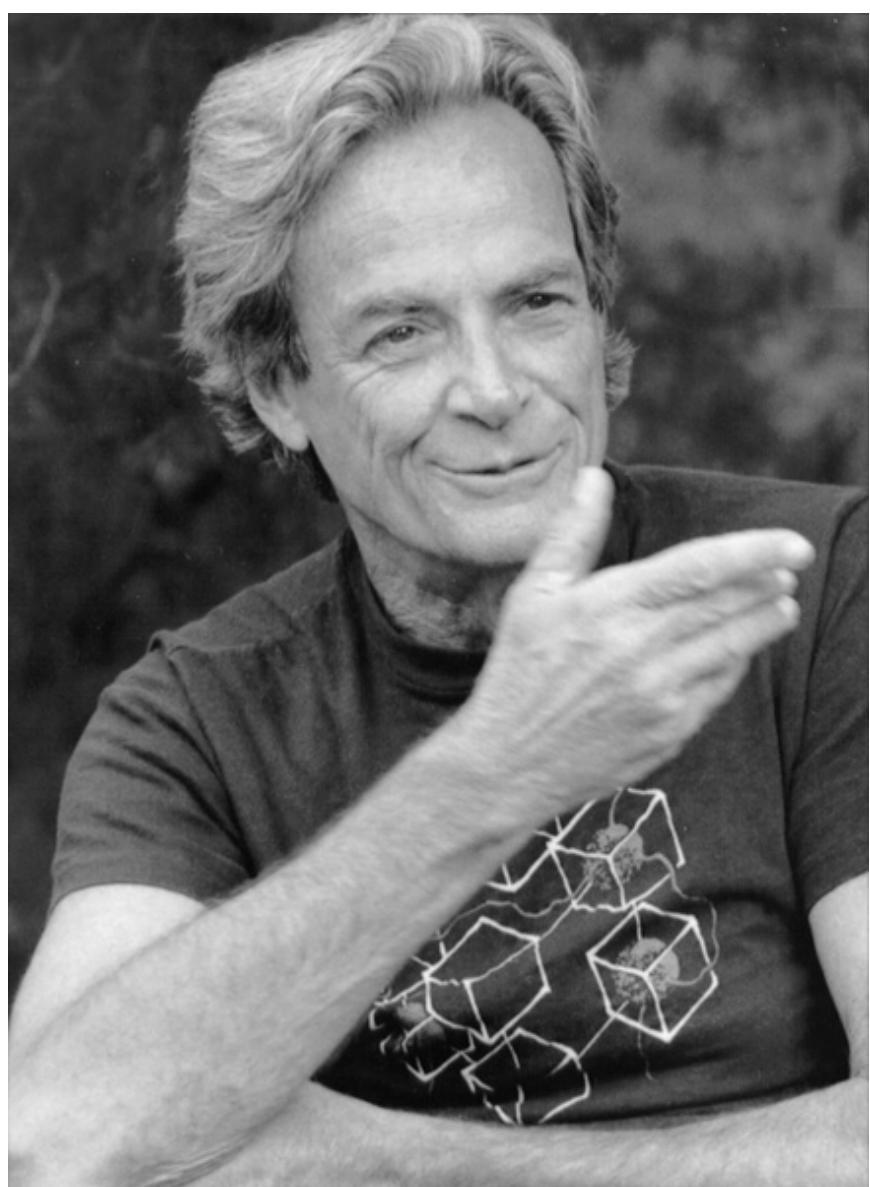
Generative learning



$$y = f(x)$$

or $p(y | x)$

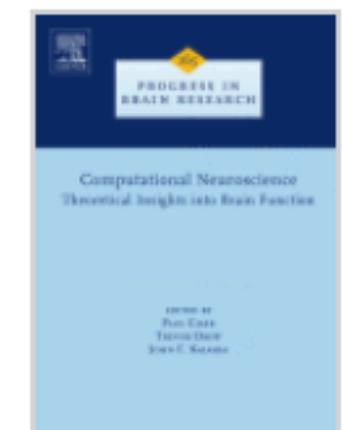
$$p(x, y)$$



Progress in Brain Research

Volume 165, 2007, Pages 535–547

Computational Neuroscience: Theoretical Insights into Brain Function



To recognize shapes, first learn to generate images

Geoffrey E. Hinton

Department of Computer Science, University of Toronto, 10 Kings College Road, Toronto, M5S 3G4
Canada

ChatGPT: Optimizing Language Models for Dialogue
November 30, 2022 — Announcements, Research

DALL·E API Now Available in Public Beta
November 3, 2022 — Announcements, API

DALL·E Now Available Without Waitlist
September 28, 2022 — Announcements

Introducing Whisper
September 21, 2022 — Research

DALL·E: Introducing Outpainting
August 31, 2022 — Announcements

Our Approach to Alignment Research
August 24, 2022 — Research

New and Improved Content Moderation Tooling
August 10, 2022 — Announcements

DALL·E Now Available in Beta
July 20, 2022 — Announcements

OpenAI Technical Goals
June 20, 2016 — Announcements

Generative Models
June 16, 2016 — Research, Milestones

Team Update
May 25, 2016 — Announcements

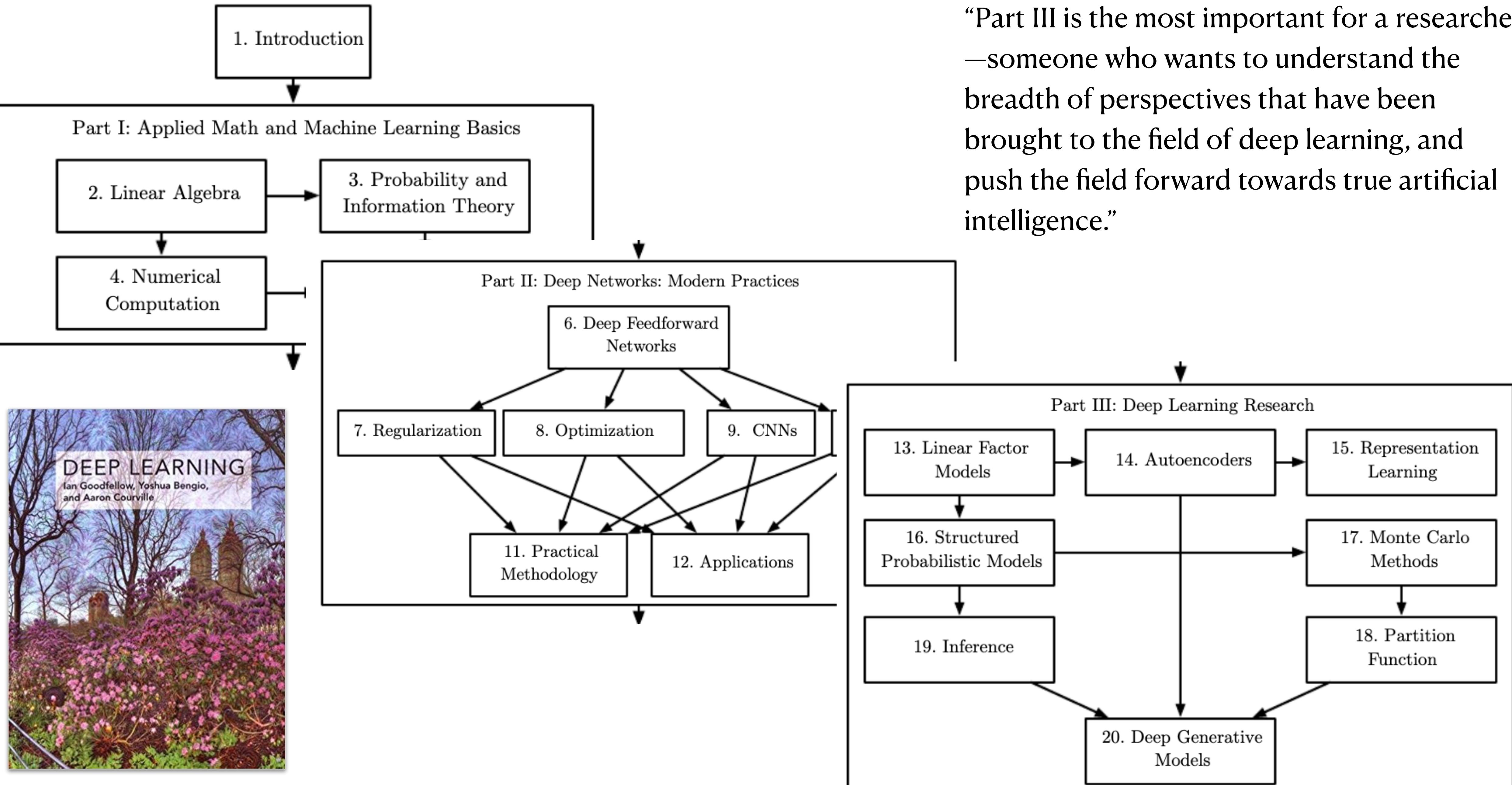
OpenAI Gym Beta
April 27, 2016 — Research

Welcome, Pieter and Shivon!
April 26, 2016 — Announcements

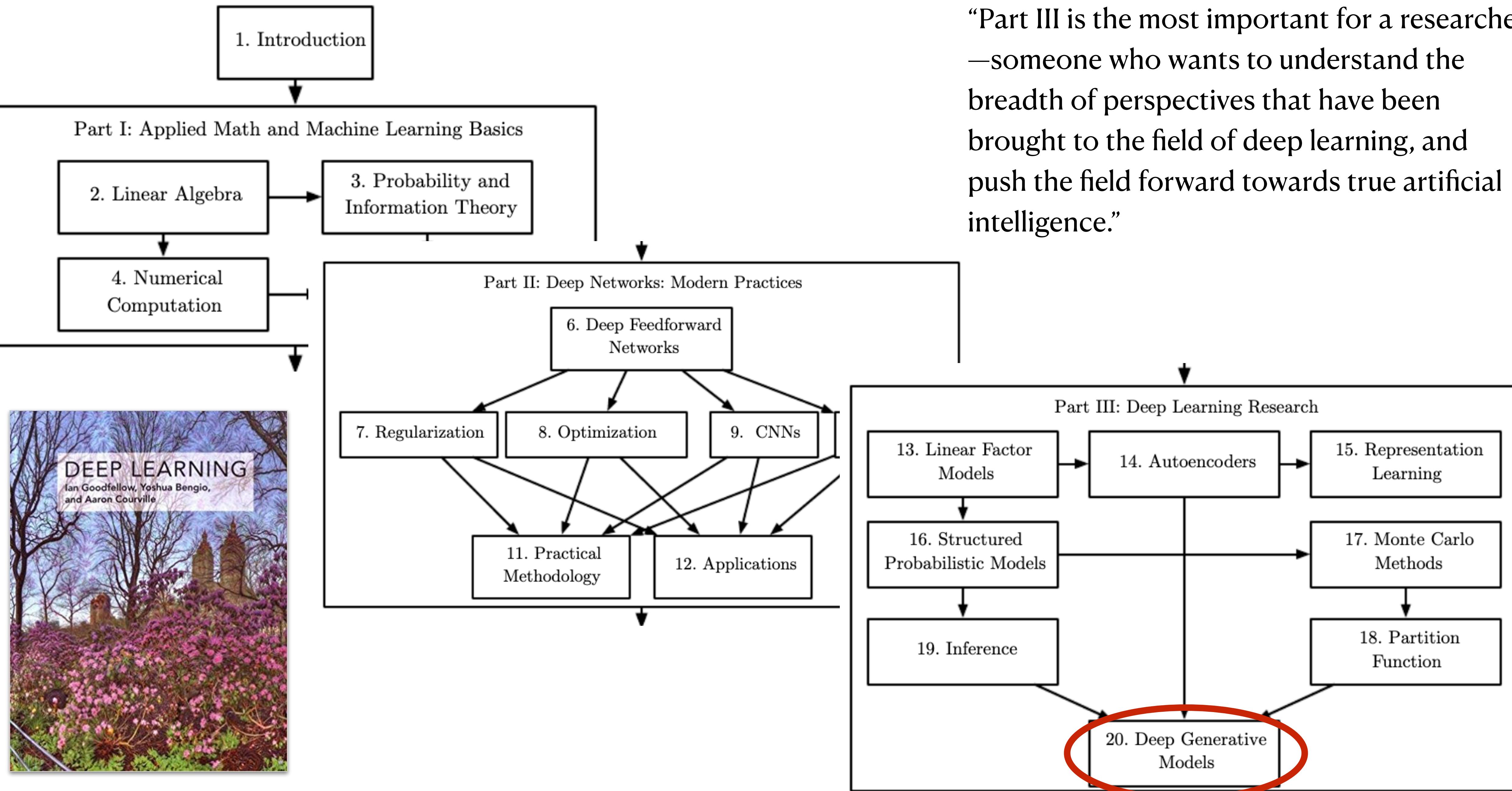
Team++
March 31, 2016 — Announcements

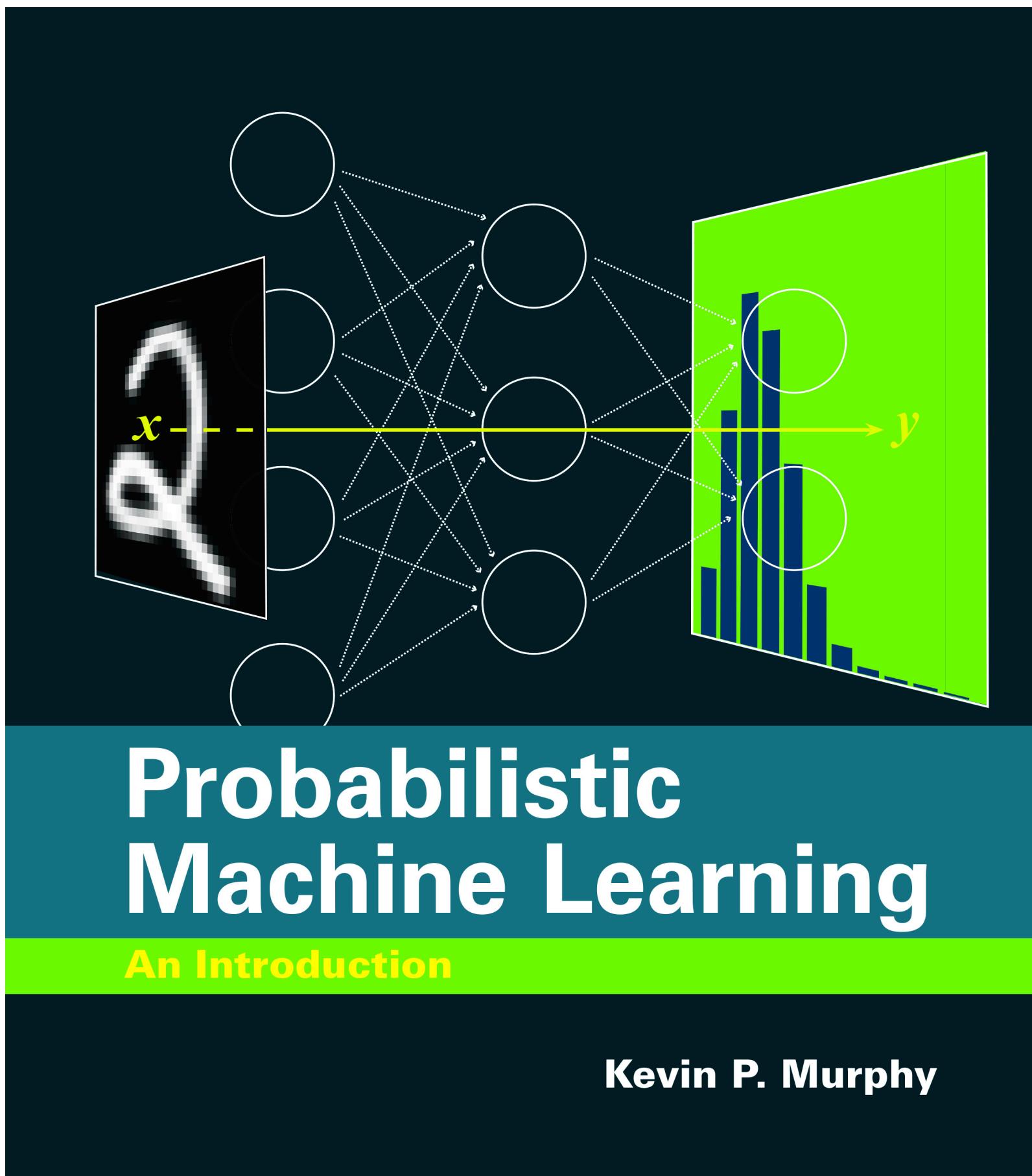
Introducing OpenAI
December 11, 2015 — Announcements

<https://openai.com/blog/>

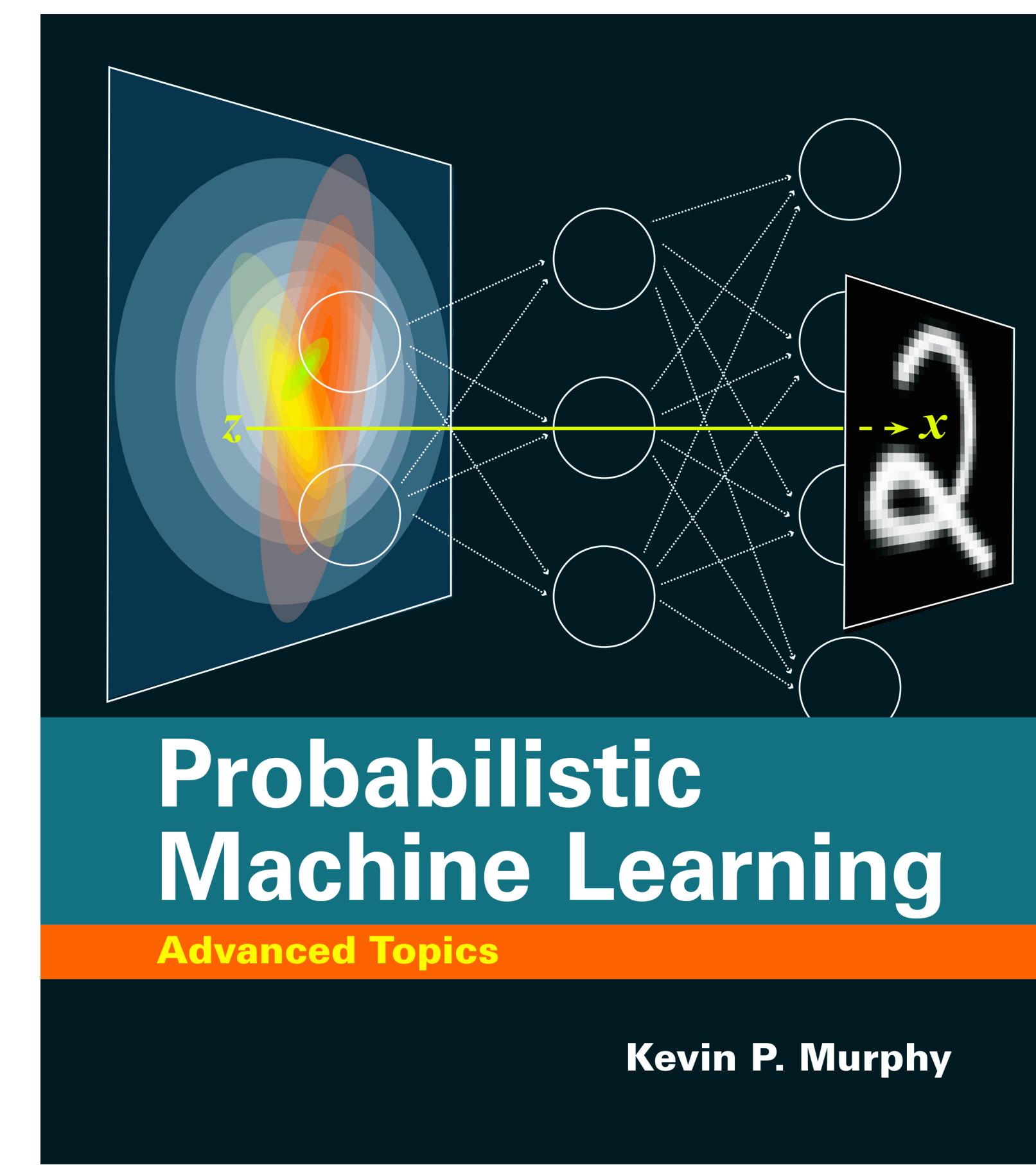


“Part III is the most important for a researcher—someone who wants to understand the breadth of perspectives that have been brought to the field of deep learning, and push the field forward towards true artificial intelligence.”



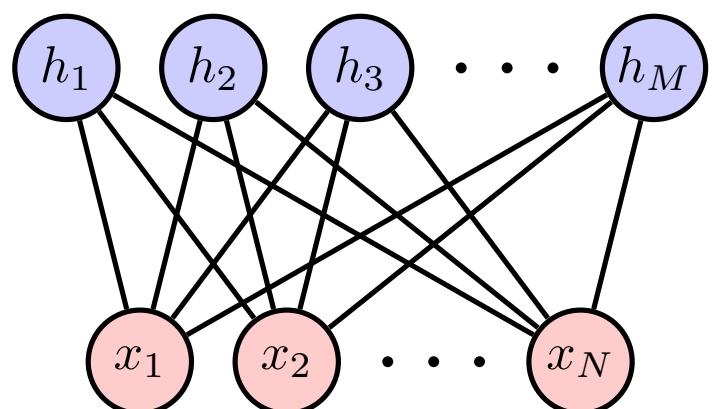


2022 (855 pages)



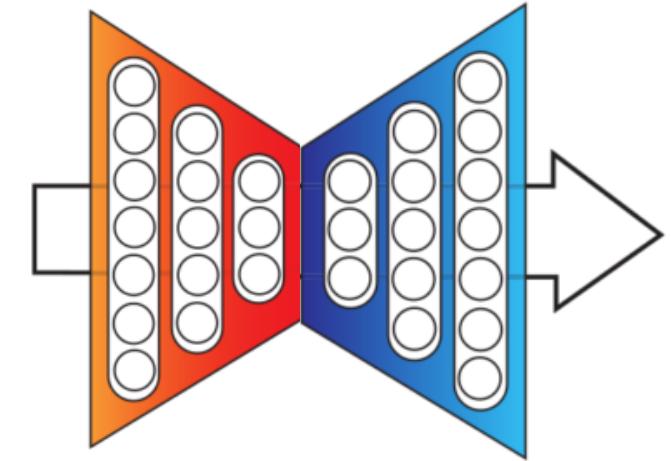
2023 (1352 +332 pages)

<https://probml.github.io/pml-book/>



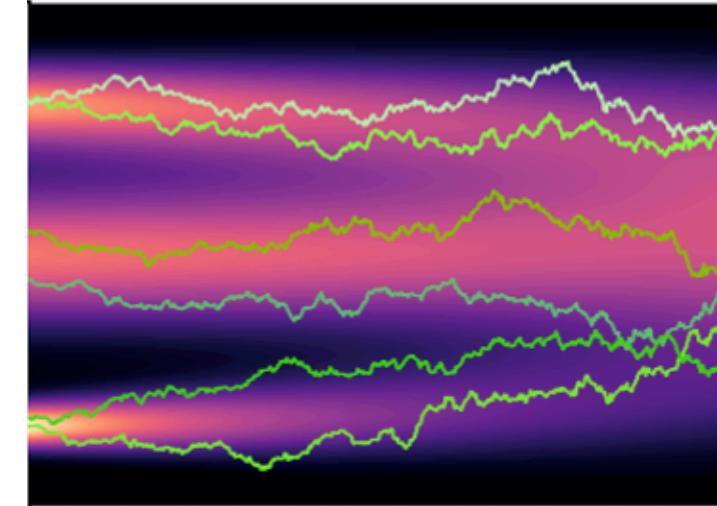
Boltzmann
Machine

1985



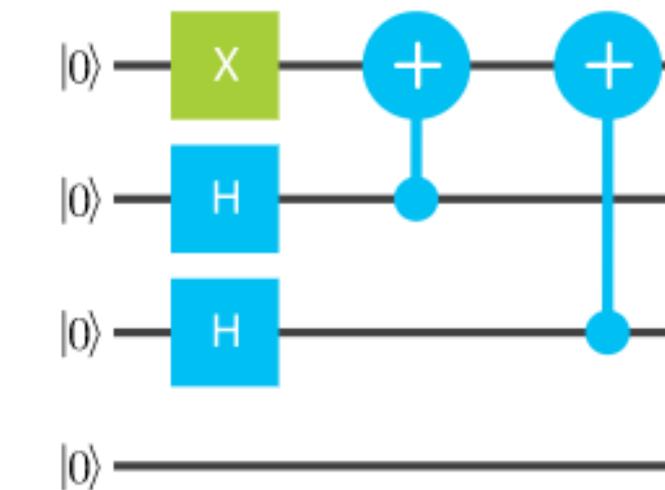
Variational
Autoencoder

2013



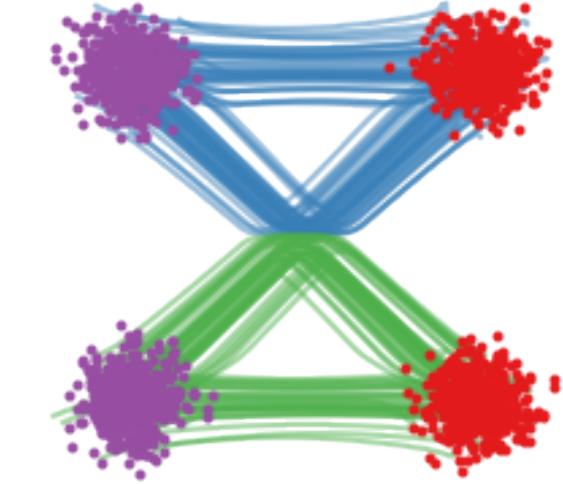
Diffusion
Model

2015



Born
Machine

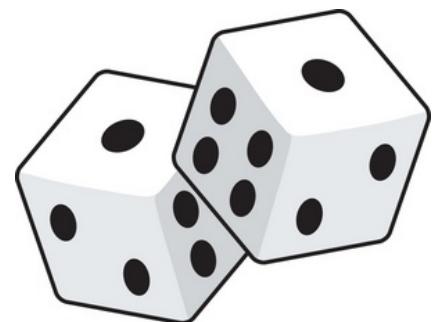
2017



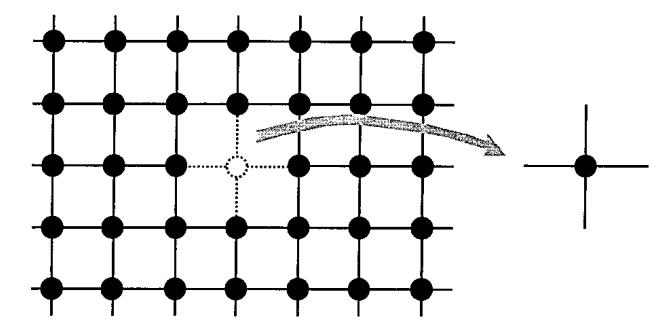
Flow
Matching

2022

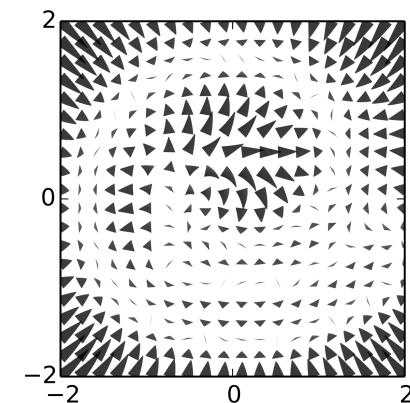
Monte Carlo
Ising model



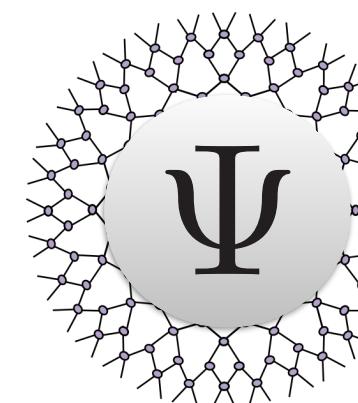
Variational
mean field



Nonequilibrium
thermodynamics



Tensor networks
Quantum circuits



Fluid optimal
transportation

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t)v] = 0$$

- ① **Leverage the power of modern generative models for physics**
- ② **Statistical, quantum, fluid, ... physics insights into generative models**

<https://future.com/how-to-build-gpt-3-for-science/>

How to Build a GPT-3 for Science

(scientific literature and data)

G_{enerative} P_{retrained} T_{ransformer}

text $\sim p(\text{text} \mid \text{prompt})$

Josh Nicholson

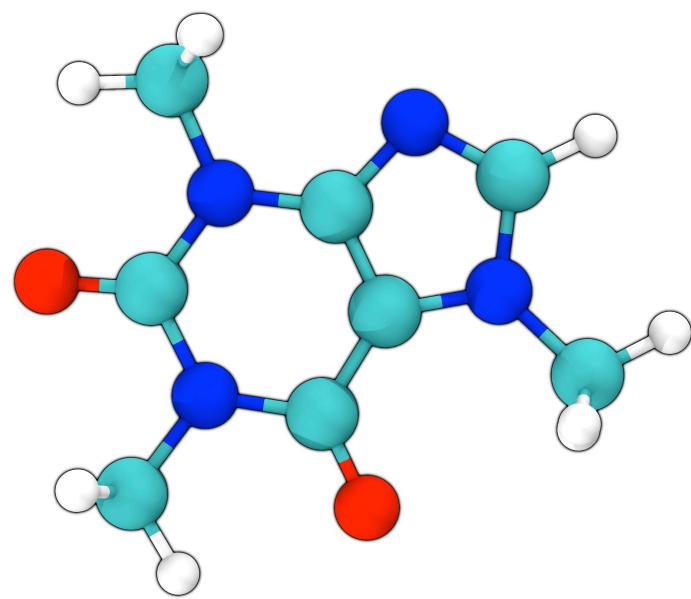
Posted August 18, 2022

Galactica, ChemGPT, MaterBERT,
ChemCrow, MatChat...

You may ask (prompts):

- “Tell me why this hypothesis is wrong”
- “Tell me why my treatment idea won’t work”
- “Generate a new treatment idea”
- “What evidence is there to support social policy X?”
- “Who has published the most reliable research in this field?”
- “Write me a scientific paper based on my data”

Language = anything you can tokenize



||

“CN1C=NC2=C1C(=O)N(C(=O)N2C)C”

Simplified Molecular-Input Line-Entry System (SMILES)

Meta AI, Galactica: A Large Language Model for Science, 2211.09085

Modality	Sequence
Text	Abell 370 is a cluster...
LAT _E X	$r_{\{s\}} = \frac{2GM}{c^2}$
Code	class Transformer(nn.Module)
SMILES	C(C(=O)O)N
AA Sequence	MIRLGAPQTL..
DNA Sequence	CGGTACCCTC..

<https://whitead.github.io/svelte-chem-algebra/>

mol_algebra

gpt-4

mutate add sub

Prompt: [omitted prefix] Add {left} and {right} into a single molecule

author: andrew

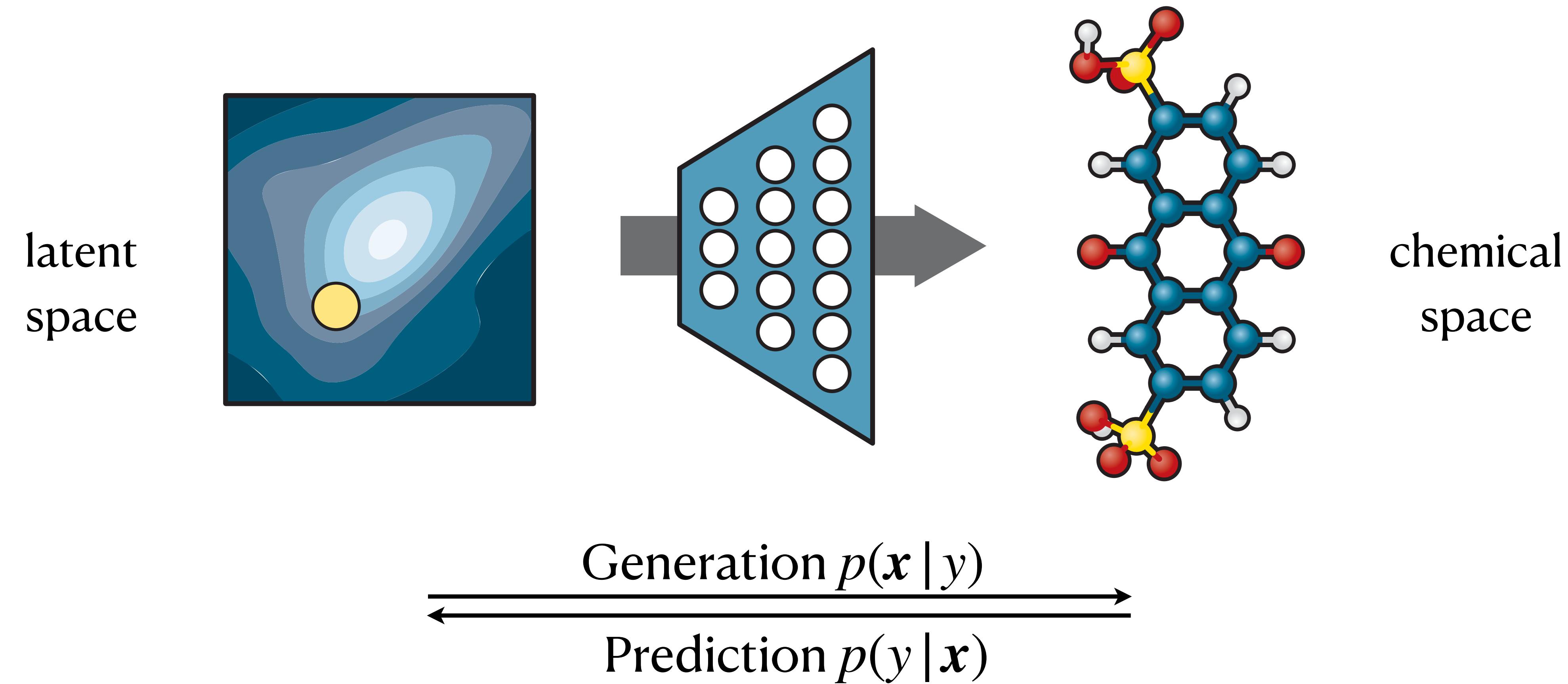
Comment | Published: 19 May 2023

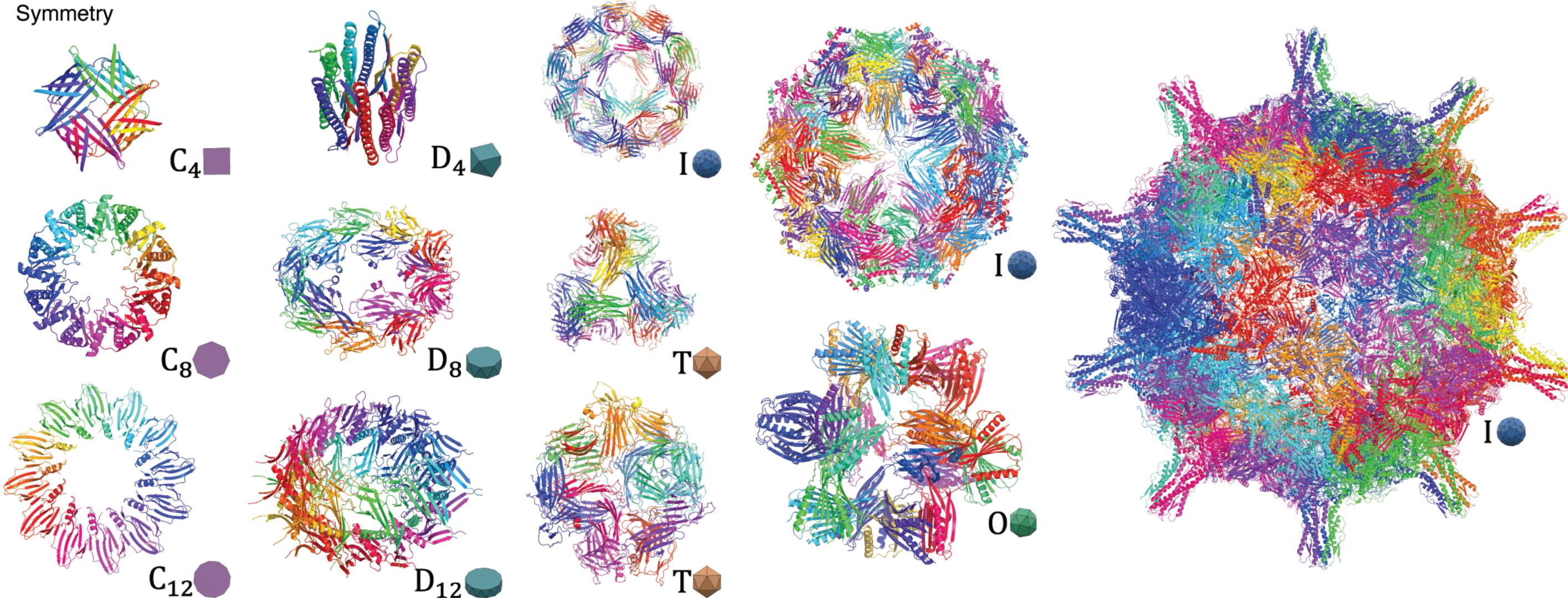
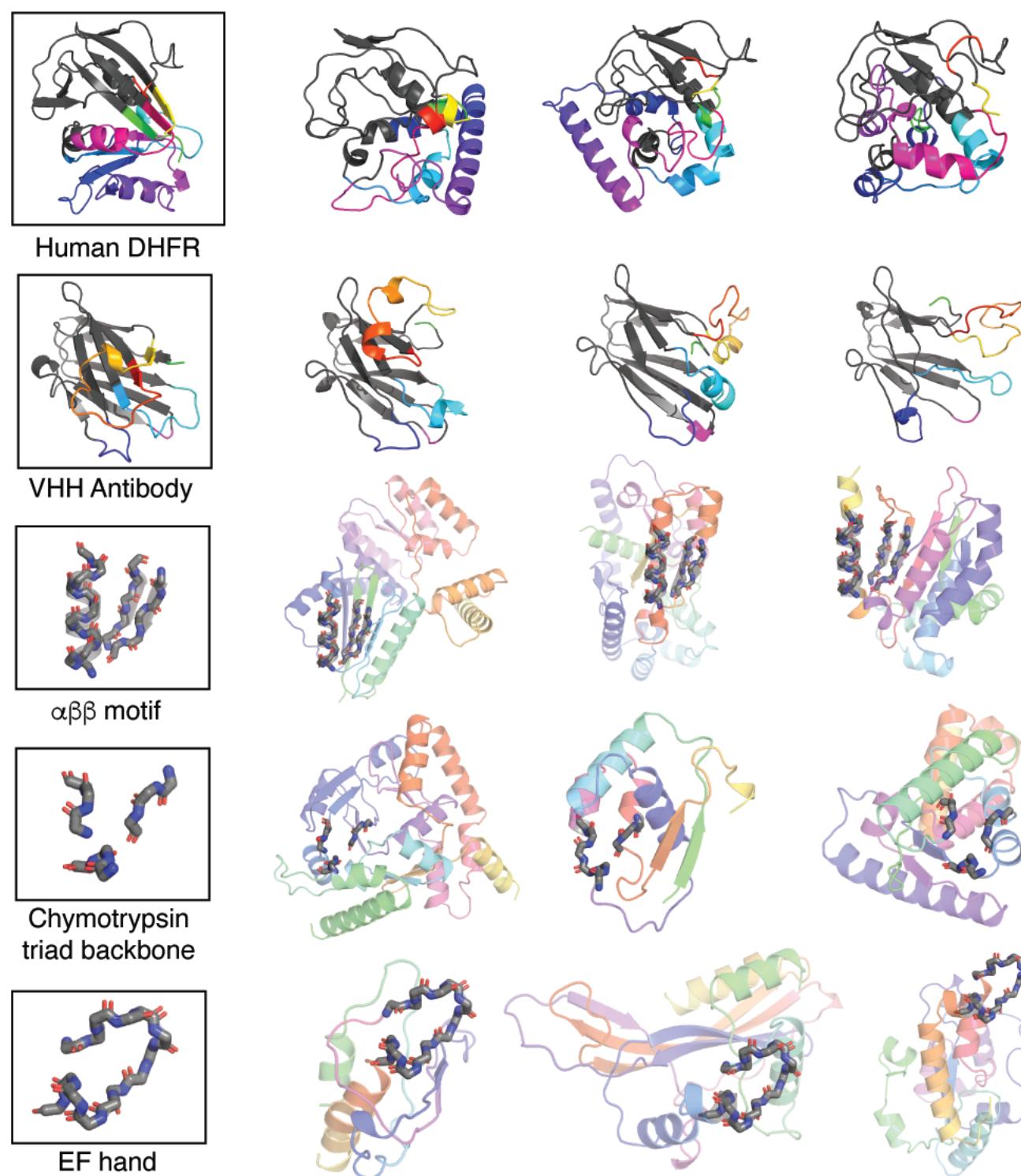
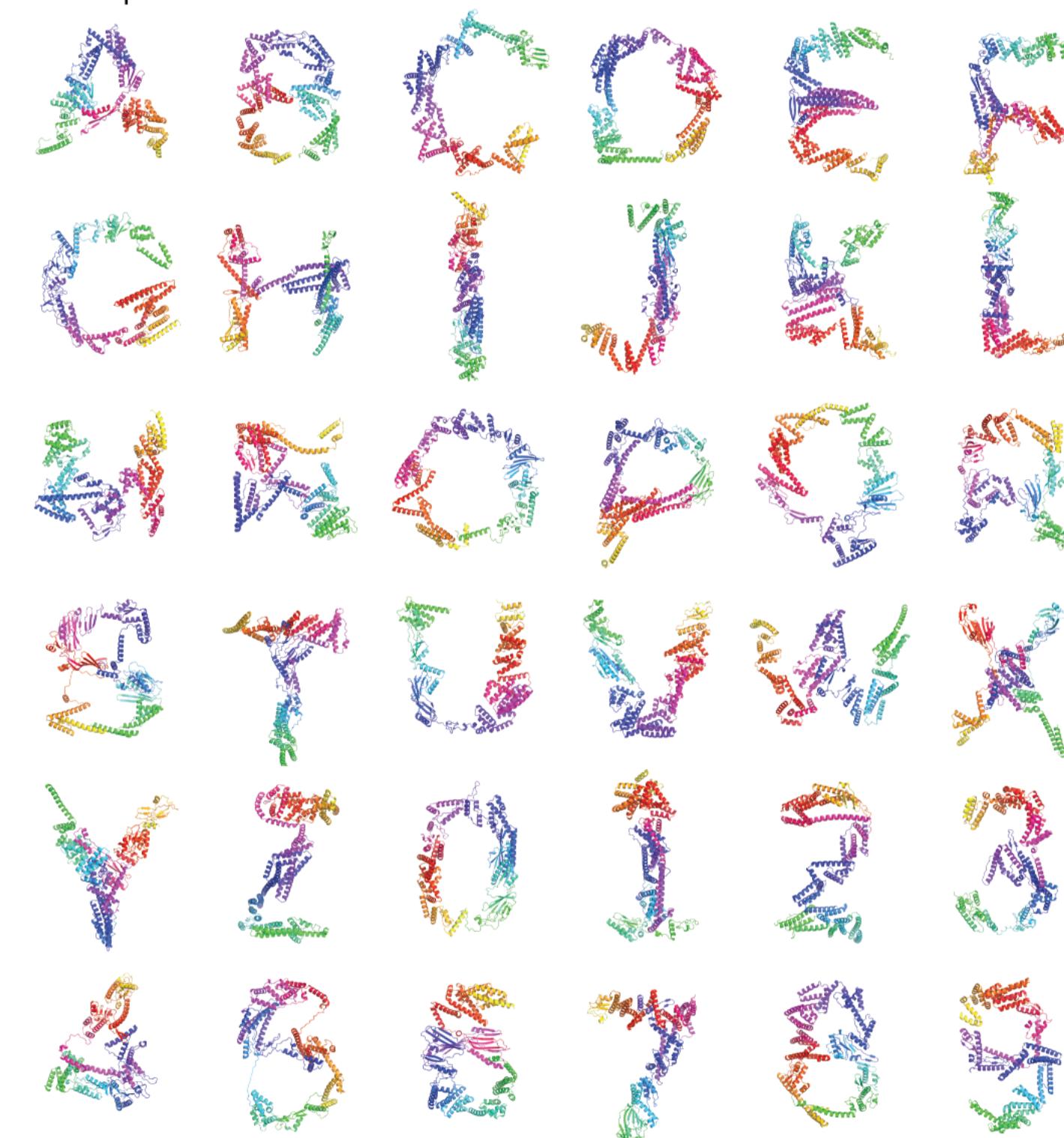
The future of chemistry is language

[Andrew D. White](#)

[Nature Reviews Chemistry](#) (2023) | [Cite this article](#)

Generative AI for matter engineering



a Symmetry**b Substructure****c Shape** $p(\text{protein} \mid \text{symmetry})$ $p(\text{protein} \mid \text{substructure})$ $p(\text{protein} \mid \text{shape})$

DeepMind

Mapping ML methods to protein problems

John Jumper

CASP15

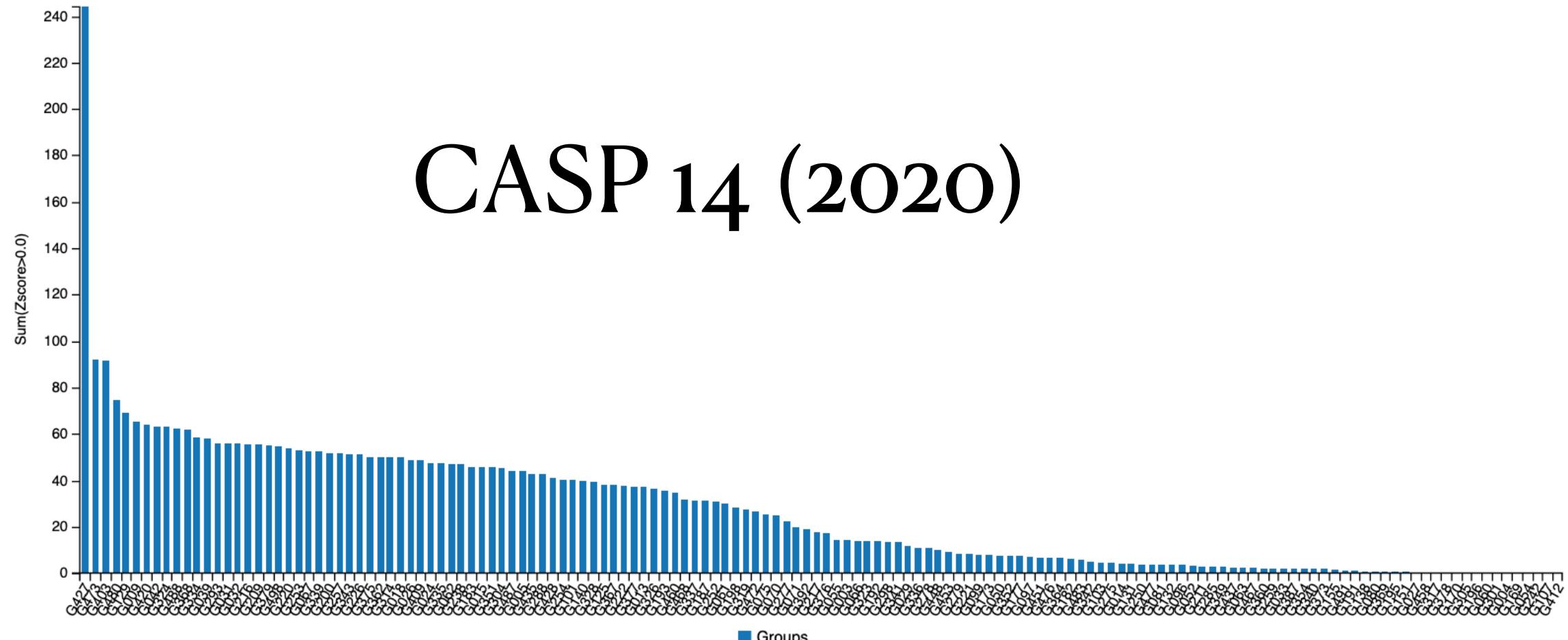


CASP 15 invited talk by John Jumper

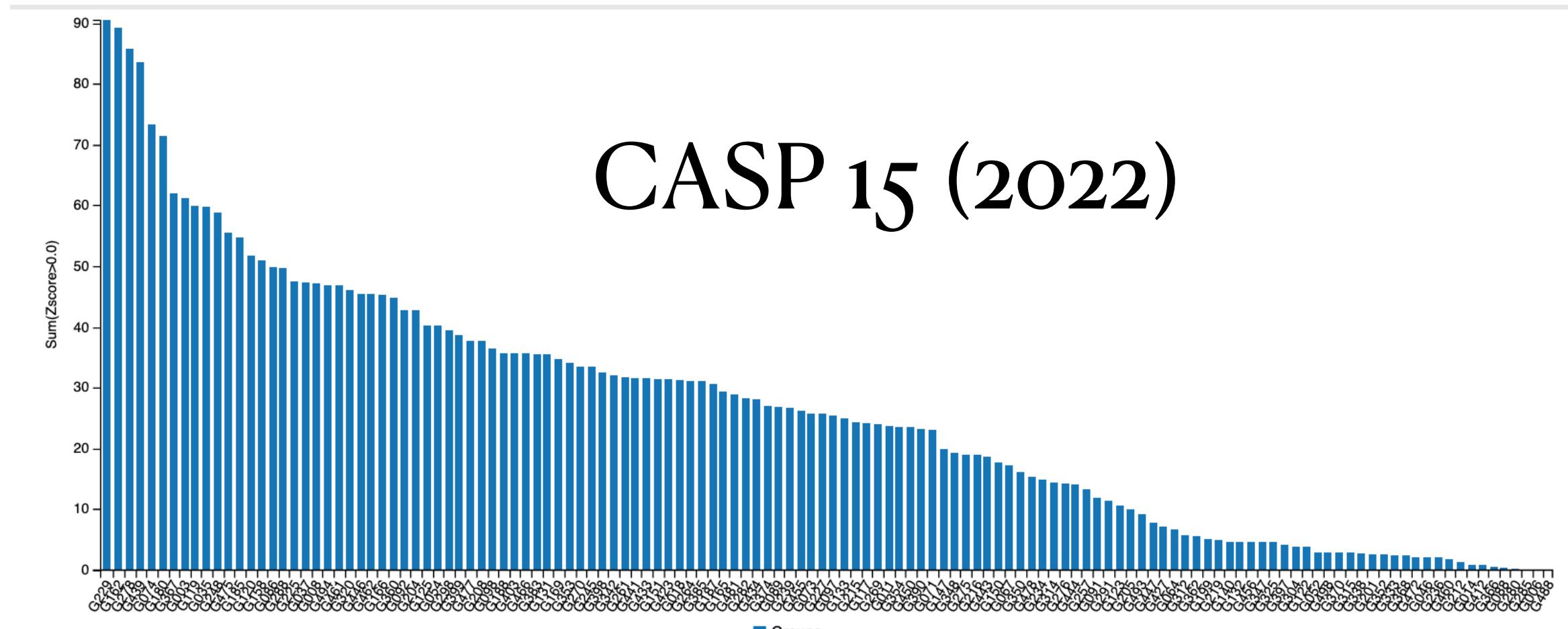
Outline

- Generative models and diffusion
- Protein language models and the scaling hypothesis
- Next problems

CASP 14 (2020)

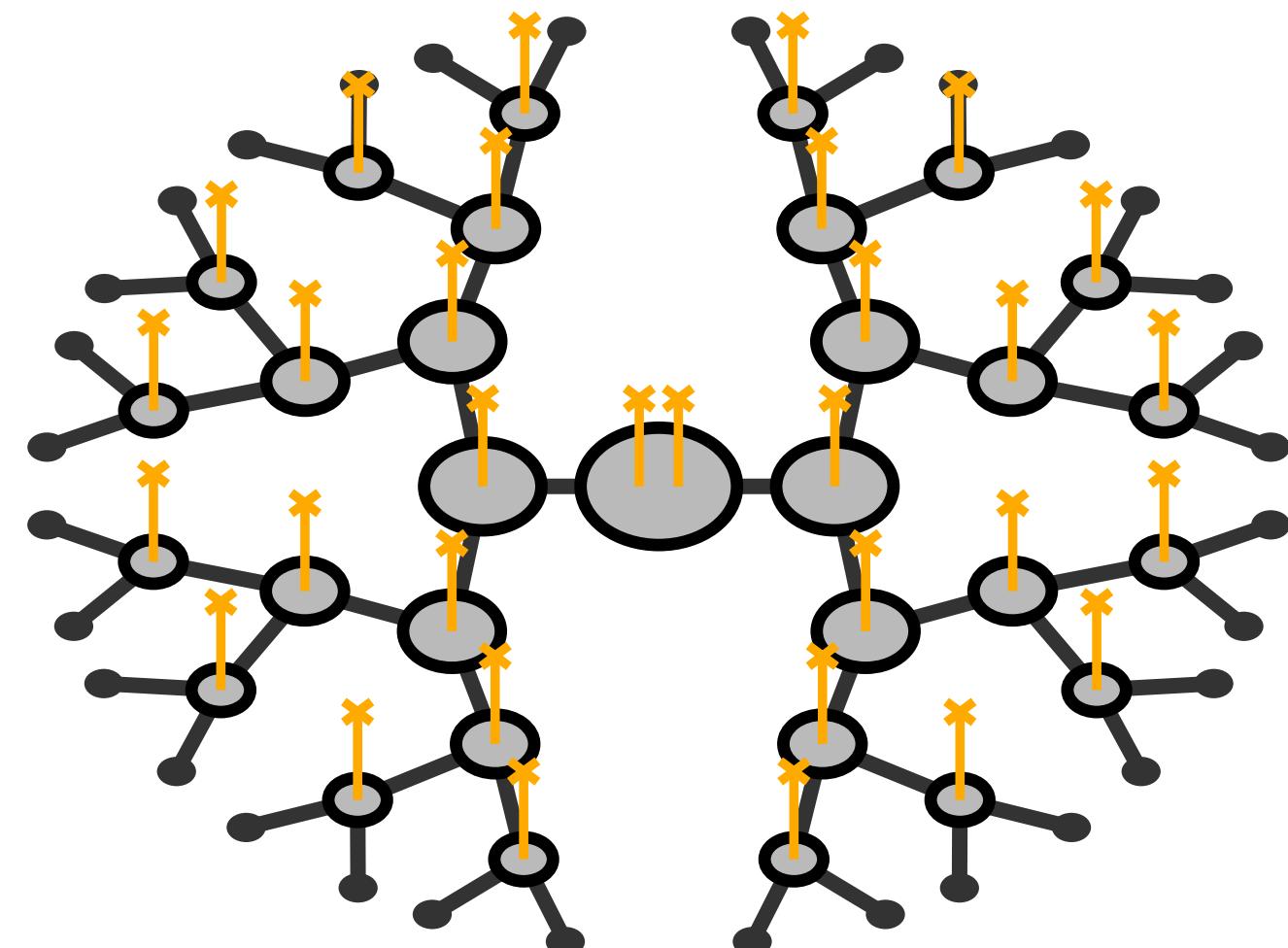


CASP 15 (2022)



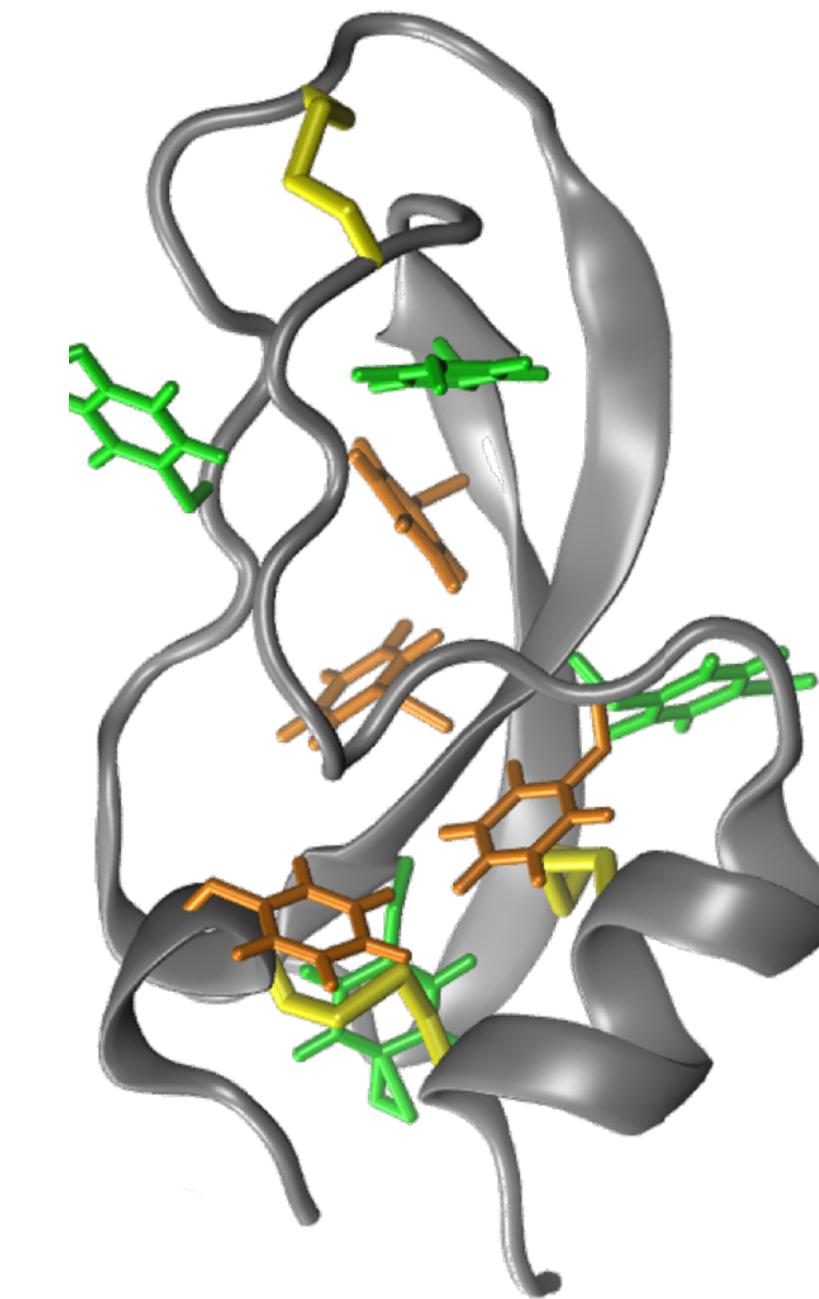
Generative AI for matter computation

Renormalization group



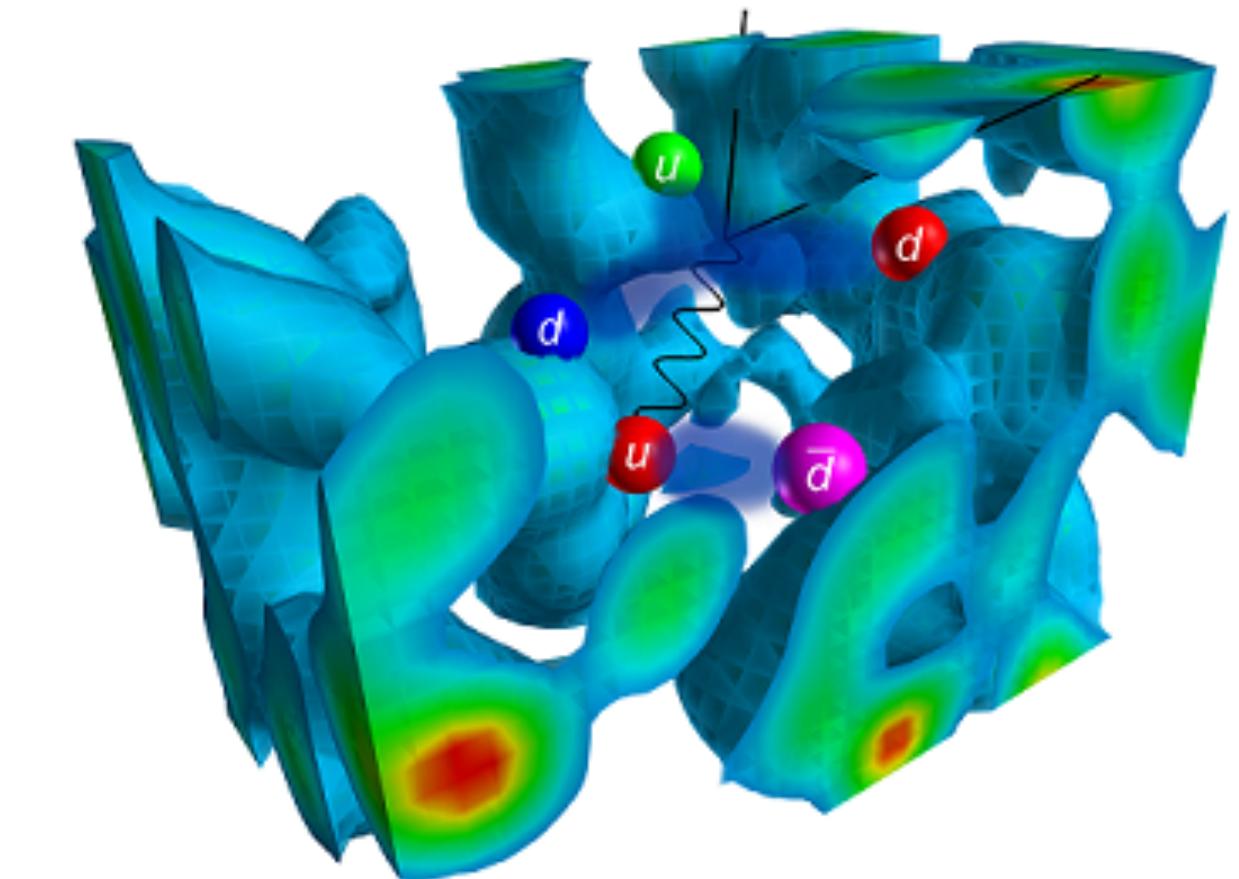
Li and LW, PRL '18
Li, Dong, Zhang, LW, PRX '20

Molecular simulation



Noe et al, Science '19
Wirnsberger et al, JCP '20

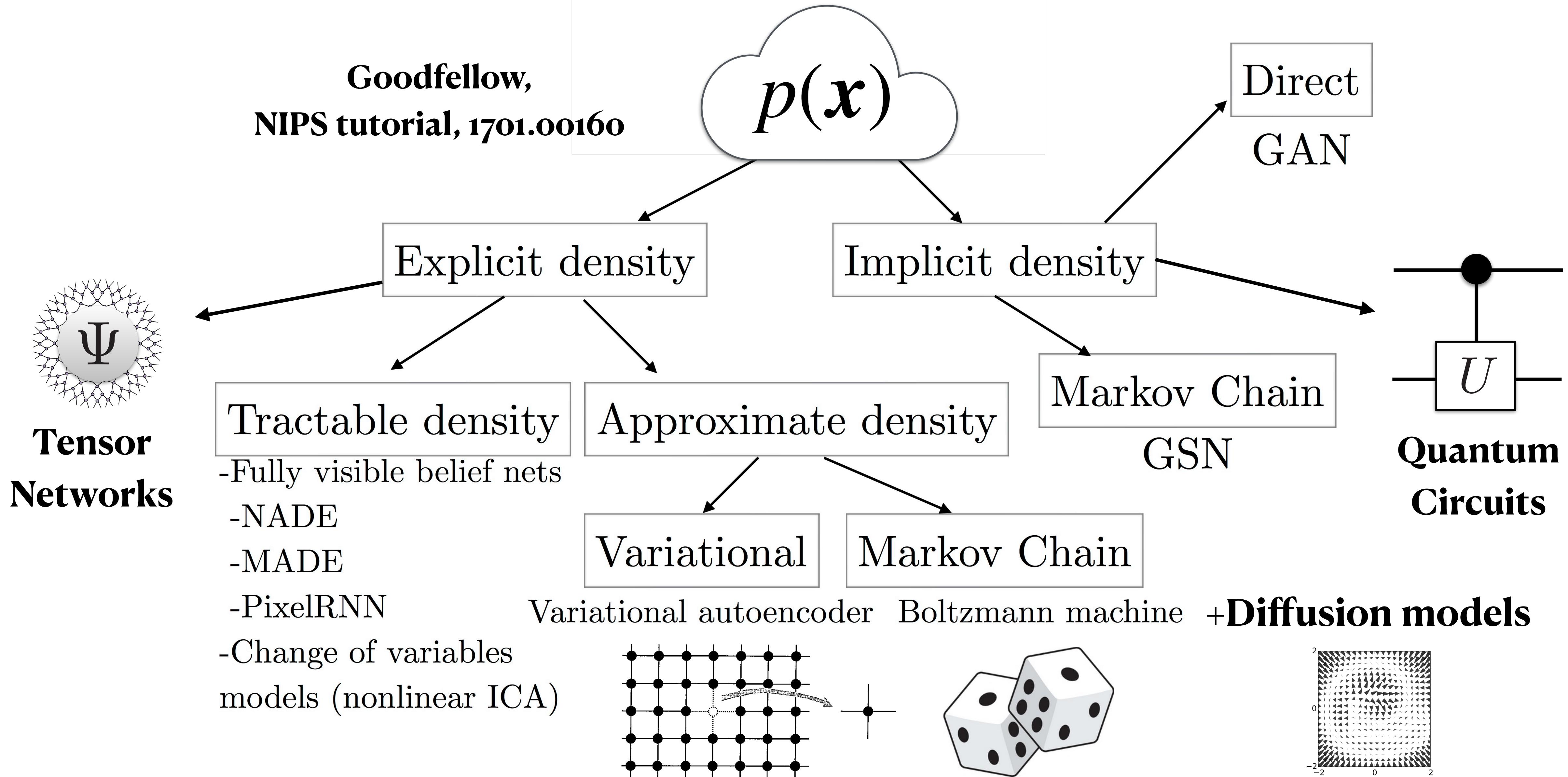
Lattice field theory



Albergo et al, PRD '19
Kanwar et al, PRL '20

These are principled calculations: quantitatively accurate,
interpretable, reliable, and generalizable even without data

Generative models and their physics genes



Probabilistic Generative Modeling

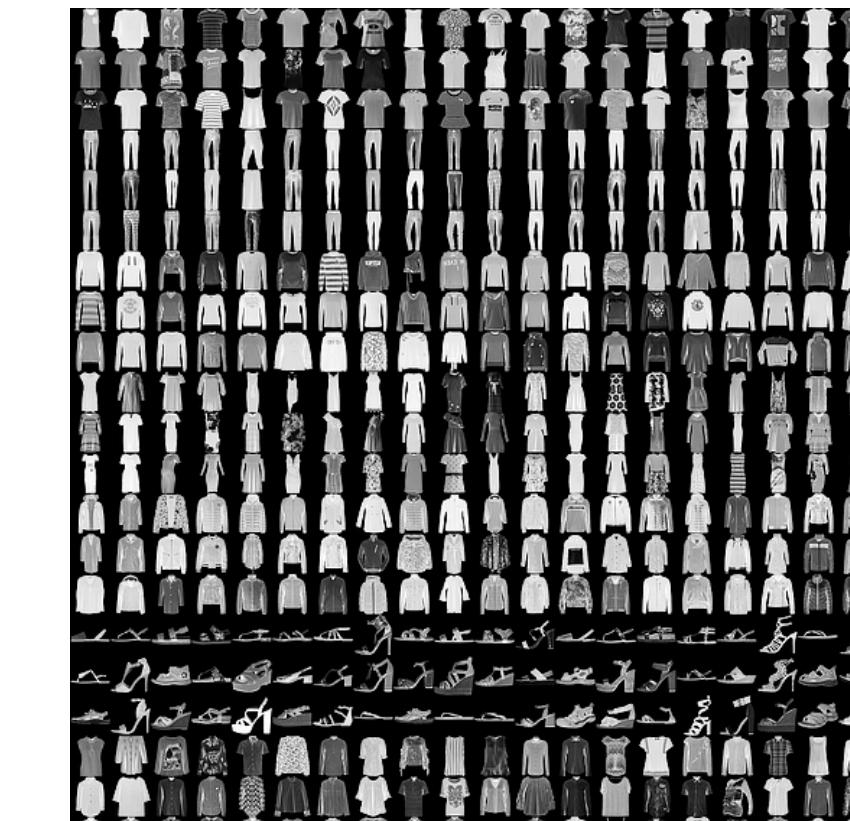
$$p(\mathbf{x})$$

How to express, learn, and sample from a
high-dimensional probability distribution ?



“random” images

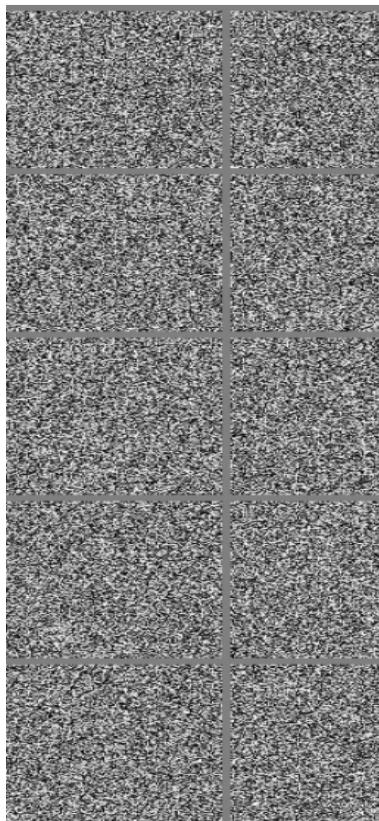
8	9	0	1	2	3	4	7	8	9	0	1	2	3	4	5	6	7	8	6
4	2	6	4	7	5	5	4	7	8	9	2	9	3	9	3	8	2	0	5
0	1	0	4	2	6	5	3	5	3	8	0	0	3	4	1	5	3	0	8
3	0	6	2	7	1	1	8	1	7	1	3	8	9	7	6	7	4	1	6
7	5	1	7	1	9	8	0	6	9	4	9	9	3	7	1	9	2	2	5
3	7	8	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	0
1	2	3	4	5	6	7	8	9	8	1	0	5	5	1	9	0	4	1	9
3	8	4	7	7	8	5	0	6	5	5	3	3	3	9	8	1	4	0	6
1	0	0	6	2	1	1	3	2	8	8	7	8	4	6	0	2	0	3	6
8	7	1	5	9	9	3	2	4	9	4	4	5	3	2	8	5	9	4	1
6	5	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7
8	9	0	1	2	3	4	5	6	7	8	9	6	4	2	6	4	7	5	5
4	7	8	9	2	9	3	9	3	8	2	0	9	8	0	5	6	0	1	0
4	2	6	5	5	5	4	3	4	1	5	3	0	8	3	0	6	2	7	1
1	8	1	7	1	3	8	5	4	2	0	9	7	6	7	4	1	6	8	4
7	5	1	2	6	7	1	9	8	0	6	9	4	9	9	6	2	3	7	1
9	2	2	5	3	7	8	0	1	2	3	4	5	6	7	8	0	1	2	3
4	5	6	7	8	0	1	2	3	4	5	6	7	8	9	2	1	2	1	3
9	9	8	5	3	7	0	7	7	5	7	9	9	4	7	0	3	4	1	4
4	7	5	8	1	4	8	4	1	8	6	4	6	3	5	7	2	5	9	



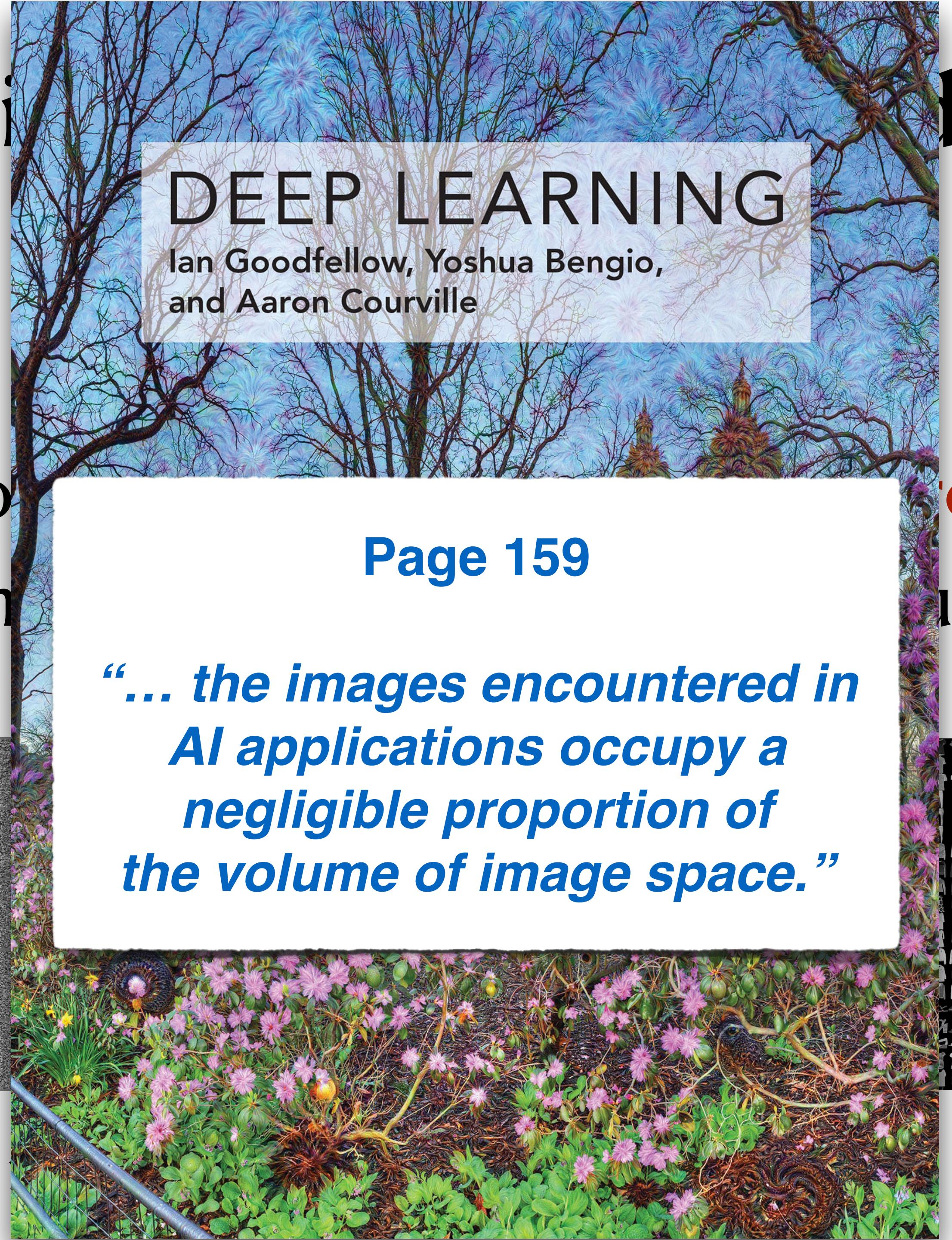
“natural” images

Probability Modeling

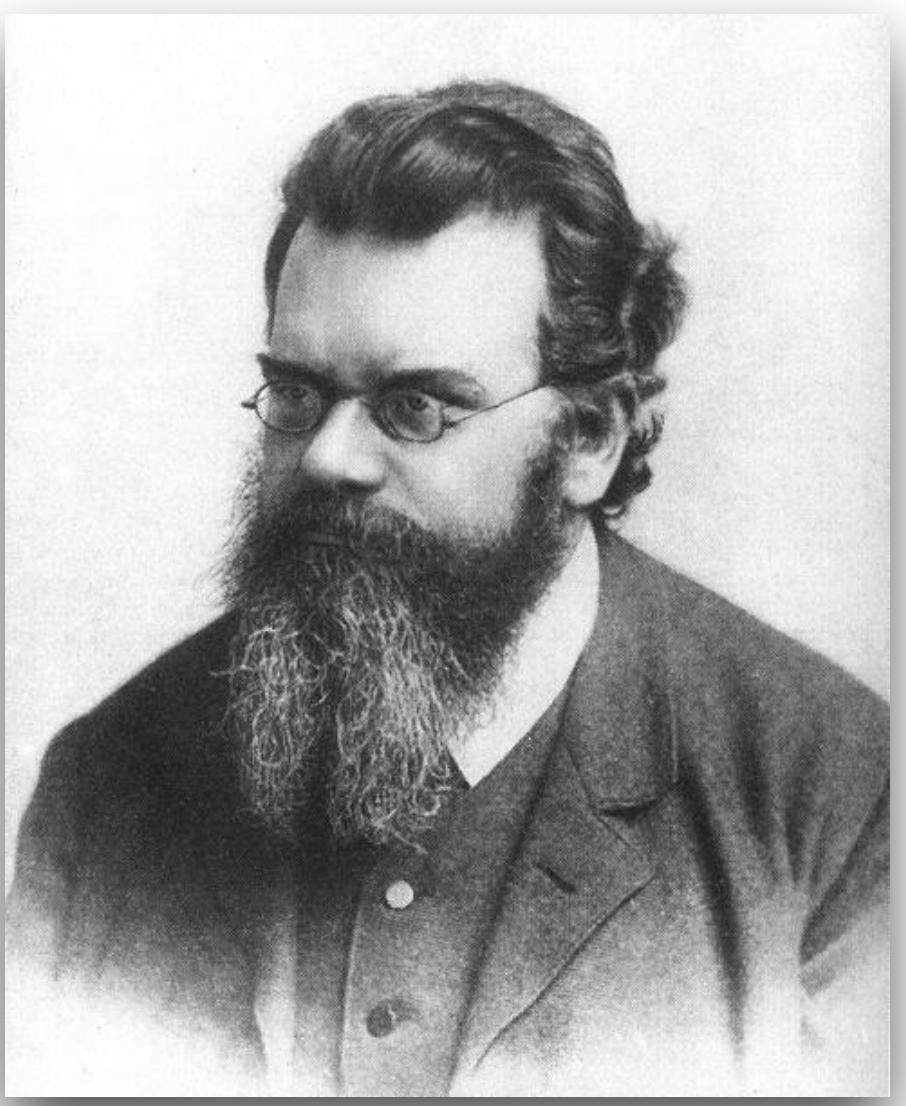
How to
high-dim



“random”



from a
dition ?



Boltzmann Machines

Ackley, Hinton, Sejnowski, Cognitive Science '85

$$p(x) = \frac{e^{-E(x)}}{Z}$$

statistical physics

“Born” Machines

Cheng, Chen, LW, Entropy '18,
Han et al, PRX 18', Liu et al PRA '18

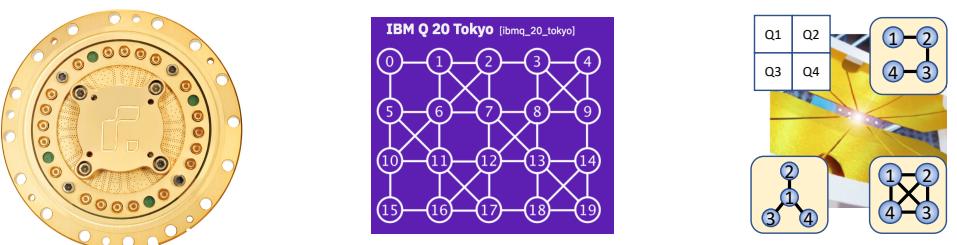
$$p(x) = |\Psi(x)|^2$$

quantum physics

Born machine: a quantum (inspired) generative model

$$p(x) = |\Psi(x)|^2$$

Quantum circuit realizations



Rigetti to build UK's first commercial quantum computer

Siddharth Venkataramakrishnan in London SEPTEMBER 2 2020

Among the first tasks for the computer is creating a “Quantum Circuit Born Machine”, said Alexei Kondratyev, managing director

IonQ and GE Research
Potential of Quantum Aggregation

June 23, 2022

COLLEGE PARK, Md., promising early results on the benefits of quantum distributions in risk man

Leveraging a Quantum Circuit Born Machine-based framework on standardized, historical indexes, IonQ and GE Research, the central innovation hub for the C

Applications of Quantum Machine Learning

Cambridge Quantum

Finance

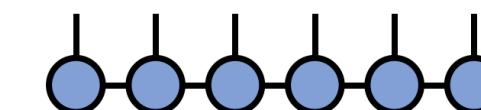
- Quantum-enhanced variational inference on hidden Markov models for time-series data
- Born Machines for foreign exchange spot return modelling
- Sampling financial data for Monte Carlo pricing using quantum GANs and Born machines

Pharmaceuticals and Healthcare

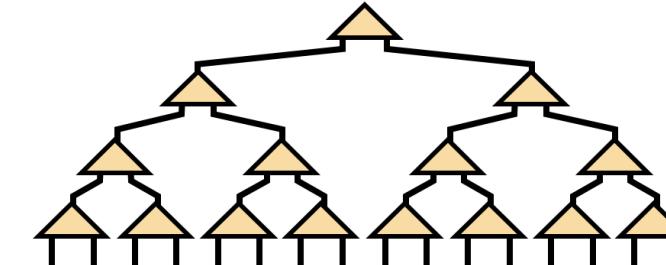
- Meta-heuristics for faster biomarker discovery in drug development based on quantum circuit Born machines
- Medical diagnosis with quantum-enhanced inference on Bayesian networks

Tensor network Born machines

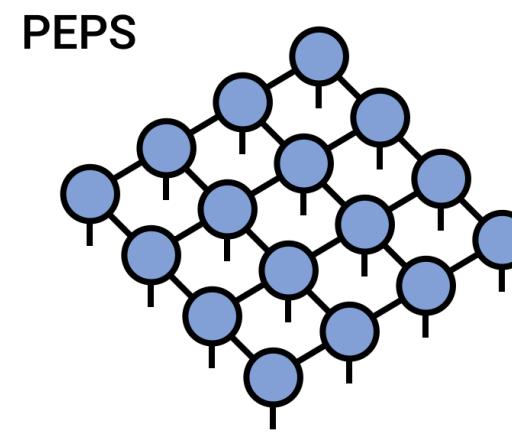
Matrix Product State / Tensor Train



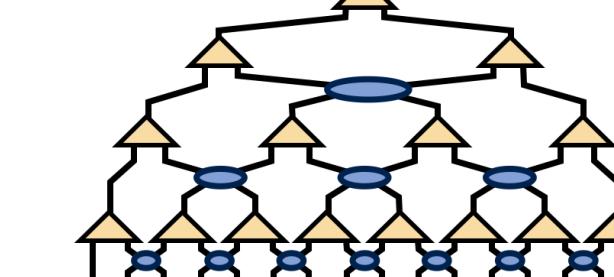
Tree Tensor Network / Hierarchical Tucker



PEPS

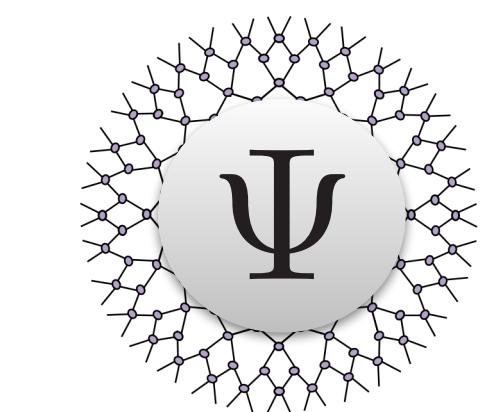


MERA



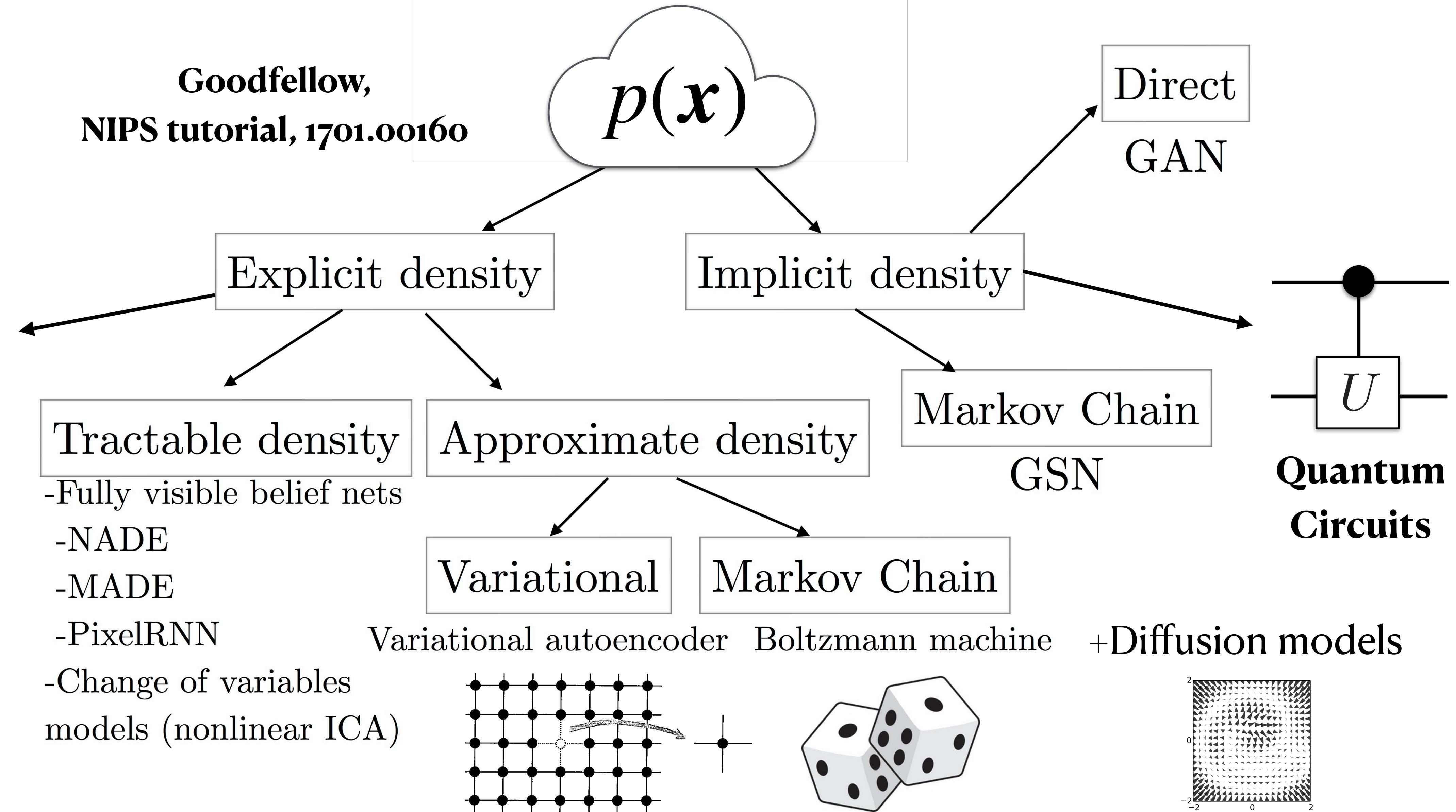
Hilbert Space
States with low entanglement

Generative models and their physics genes

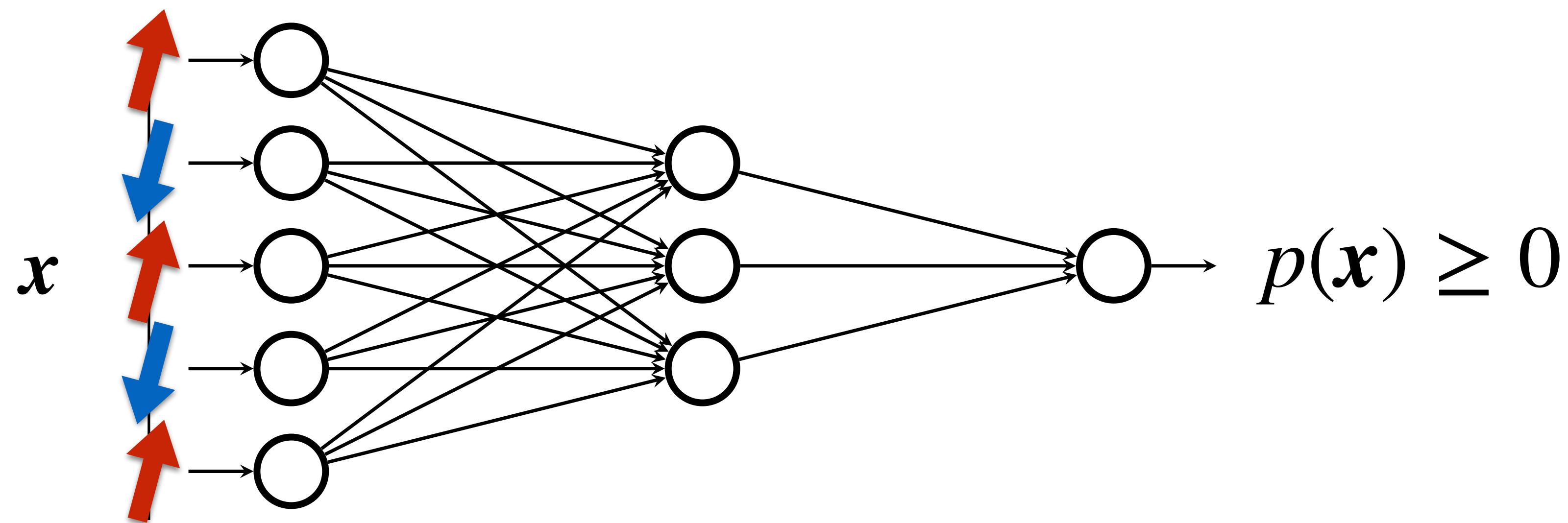


**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



So, why bother ?



Normalization ?

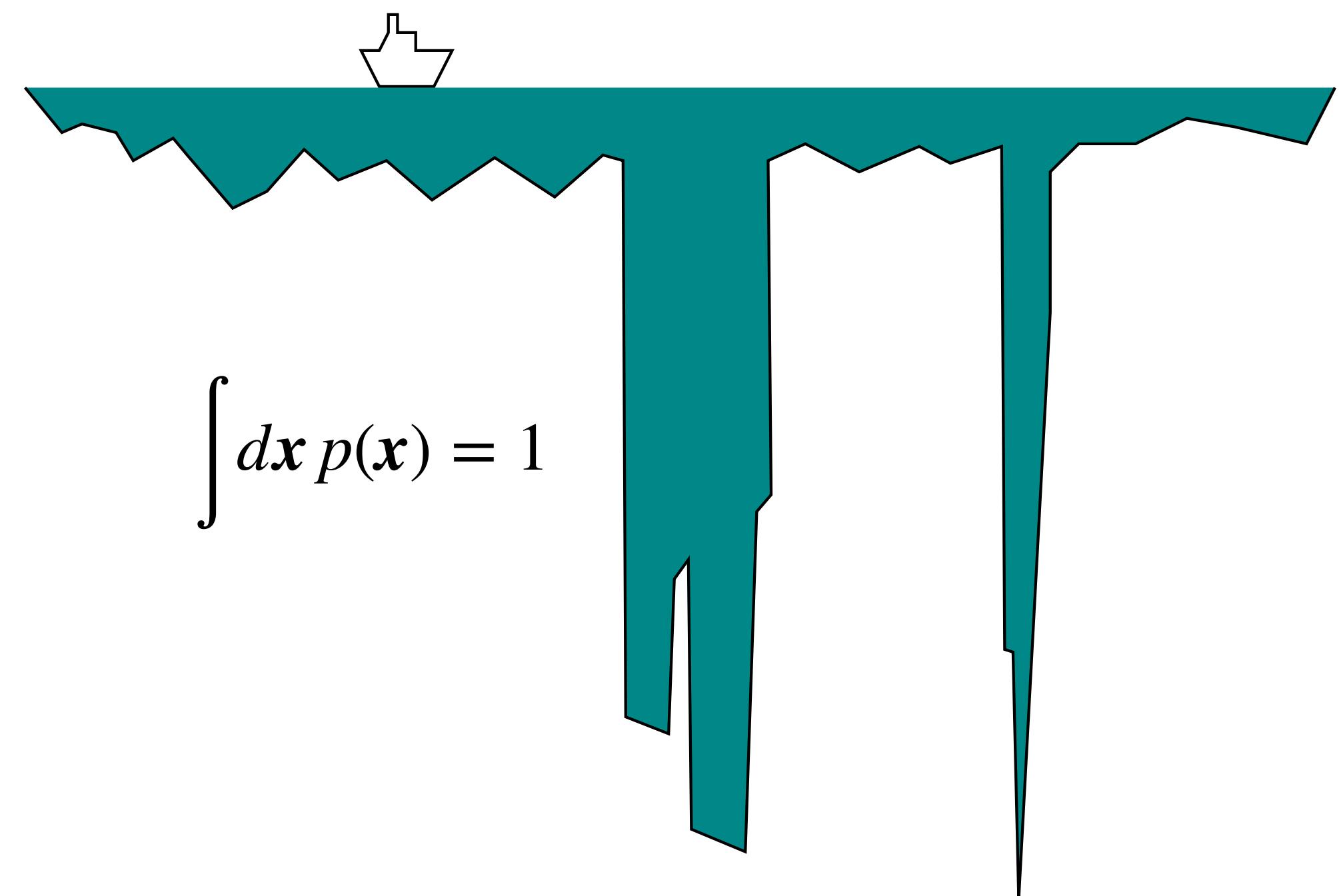
$$\int dx \, p(x)$$

Sampling ?

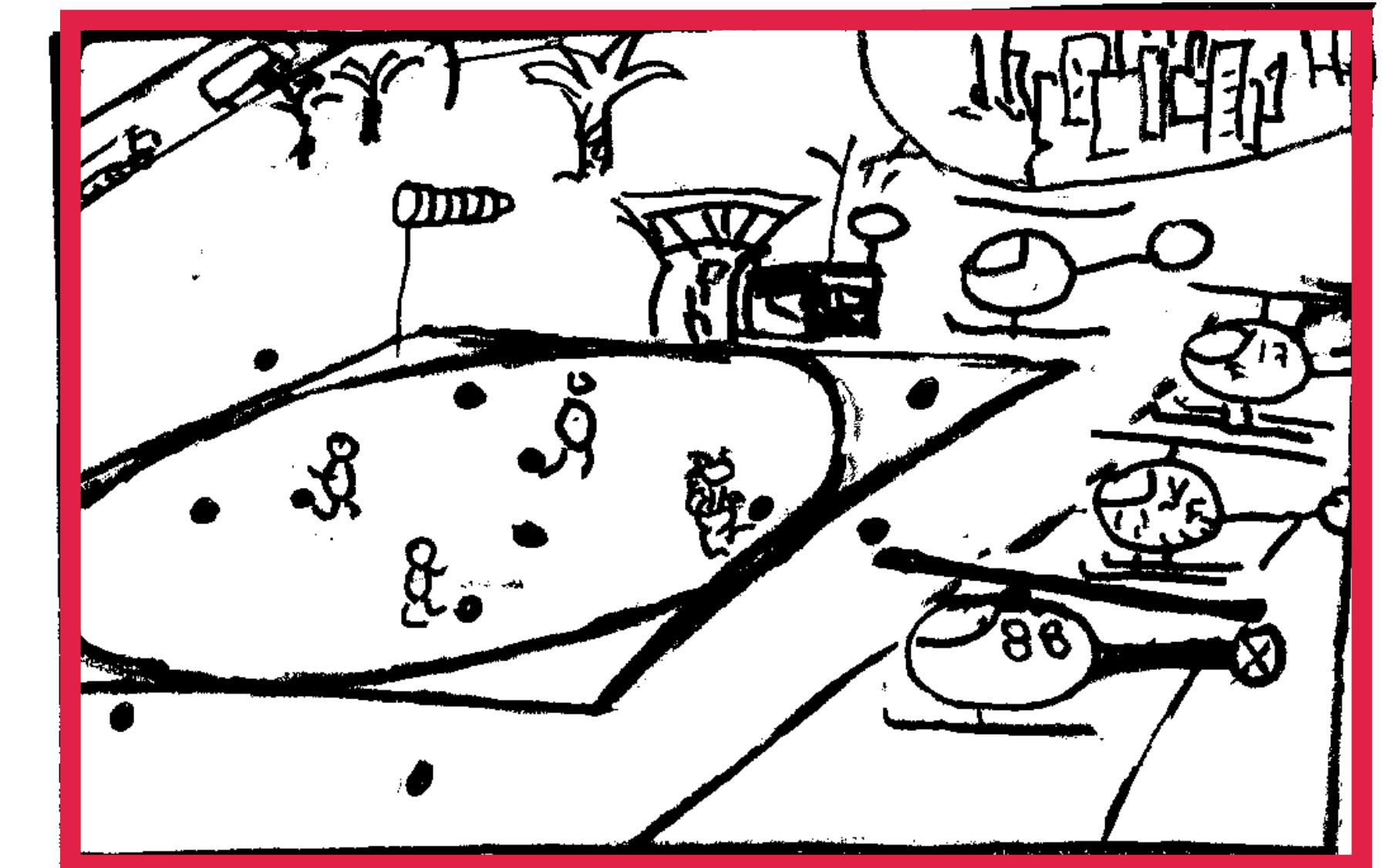
$$\mathbb{E}_{x \sim p(x)}$$

So, why bother?

Normalization



Sampling



Mackay, Information Theory, Inference, and Learning Algorithms

Krauth, Statistical Mechanics: Algorithms and Computations

We are going to see how modern generative models resolve these two issues

Generative models

Negative log-likelihood

Score function

Latent variables

Partition function

Sample diversity

Statistical physics

Energy function

Force

Collective variables/coarse
graining/renormalization group

Free energy calculation

Enhanced sampling

Two sides of the same coin

Generative modeling



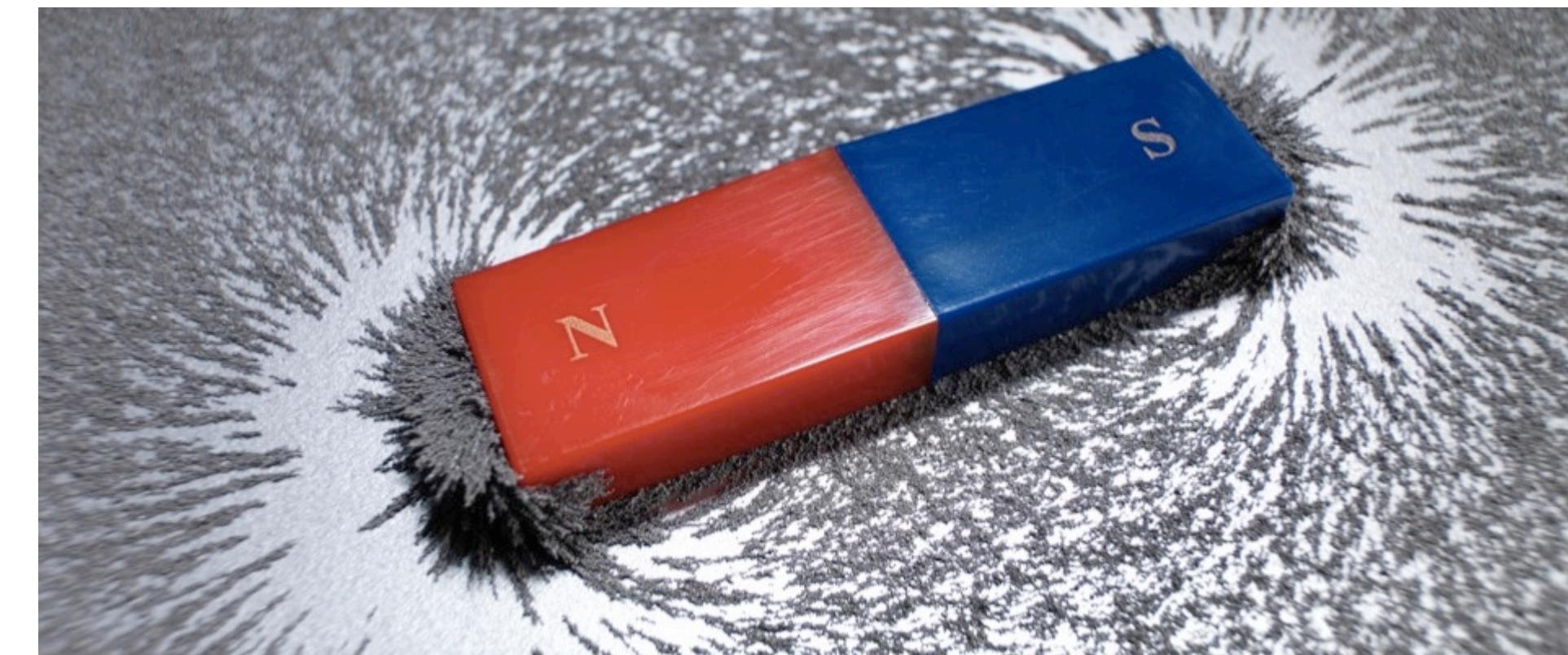
Known: samples

Unknown: generating distribution

“learn from data”

$$\mathcal{L} = - \mathbb{E}_{x \sim \text{data}} [\ln p(x)]$$

Statistical physics



Known: energy function

Unknown: samples, partition function

“learn from energy”

$$F = \mathbb{E}_{x \sim p(x)} [E(x) + k_B T \ln p(x)]$$

$$\mathbb{KL}(\text{data} \parallel p) \text{ vs } \mathbb{KL}(p \parallel e^{-E/k_B T})$$

Kullback–Leibler divergence

$$\mathbb{KL}(\pi \parallel p) \equiv \int dx \pi(x) [\ln \pi(x) - \ln p(x)]$$

$$\mathbb{KL}(\pi \parallel p) \geq 0$$

$$\mathbb{KL}(\pi \parallel p) = 0 \iff \pi(x) = p(x)$$

$$\mathbb{KL}(\pi \parallel p) \neq \mathbb{KL}(p \parallel \pi)$$

Learn from data

$$\pi(x) \propto \sum_{d \in \text{data}} \delta(x - d)$$

$$\min_{\theta} \text{KL}(\pi \| p_{\theta}) \iff \min_{\theta} \left\{ -\mathbb{E}_{x \sim \text{data}} [\ln p_{\theta}(x)] \right\}$$

target model

The lower bound is the entropy of the dataset: complete memorization

Learn from Energy

$$\pi(x) \propto e^{-E/k_B T}$$

$$\min_{\theta} \text{KL}(p_{\theta} \parallel \pi) \iff \min_{\theta} \left\{ \mathbb{E}_{x \sim p_{\theta}(x)} [E(x) + k_B T \ln p_{\theta}(x)] \right\}$$

↑ ↑

model target

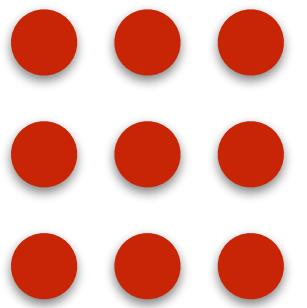
Variational free energy

The lower bound is the true free energy: exact solution

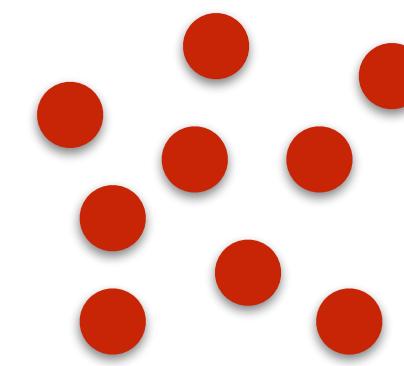
Nature tries to minimize free energy

$$F = E - TS$$

energy



entropy



F is a **cost function** of Nature

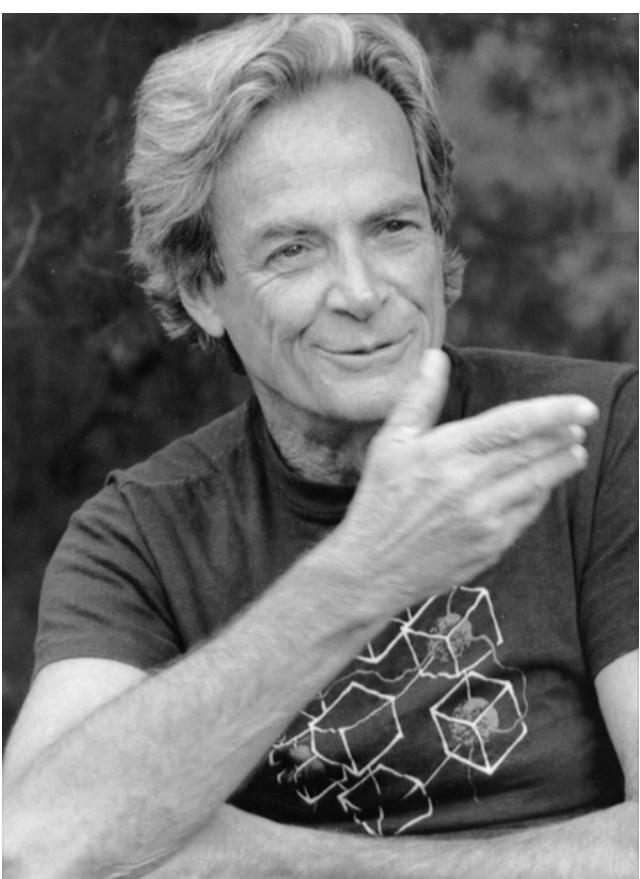
Almost the ***same*** cost function for training deep generative models

The variational free energy principle

Gibbs–Bogolyubov–Feynman

$$F[p] = \int dx p(x) [E(x) + k_B T \ln p(x)] \geq F$$

↓ ↓ ↓
variational density energy entropy 



**Difficulties in Applying the Variational
Principle to Quantum Field Theories¹**

Richard P. Feynman

Generative
models!

¹transcript of his talk in 1987

Deep variational free energy approach

Deep generative models unlock the power of
the Gibbs–Bogolyubov–Feynman–variational principle

$$F[p] = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [E(\mathbf{x}) + k_B T \ln p(\mathbf{x})]$$

↓ ↓
energy entropy 😊

Li and LW, PRL '18
Wu, LW, Zhang, PRL '19
with normalizing flow &
autoregressive models



Tractable entropy



Direct sampling



Turning a sampling problem to an optimization problem
better leverages the deep learning engine:

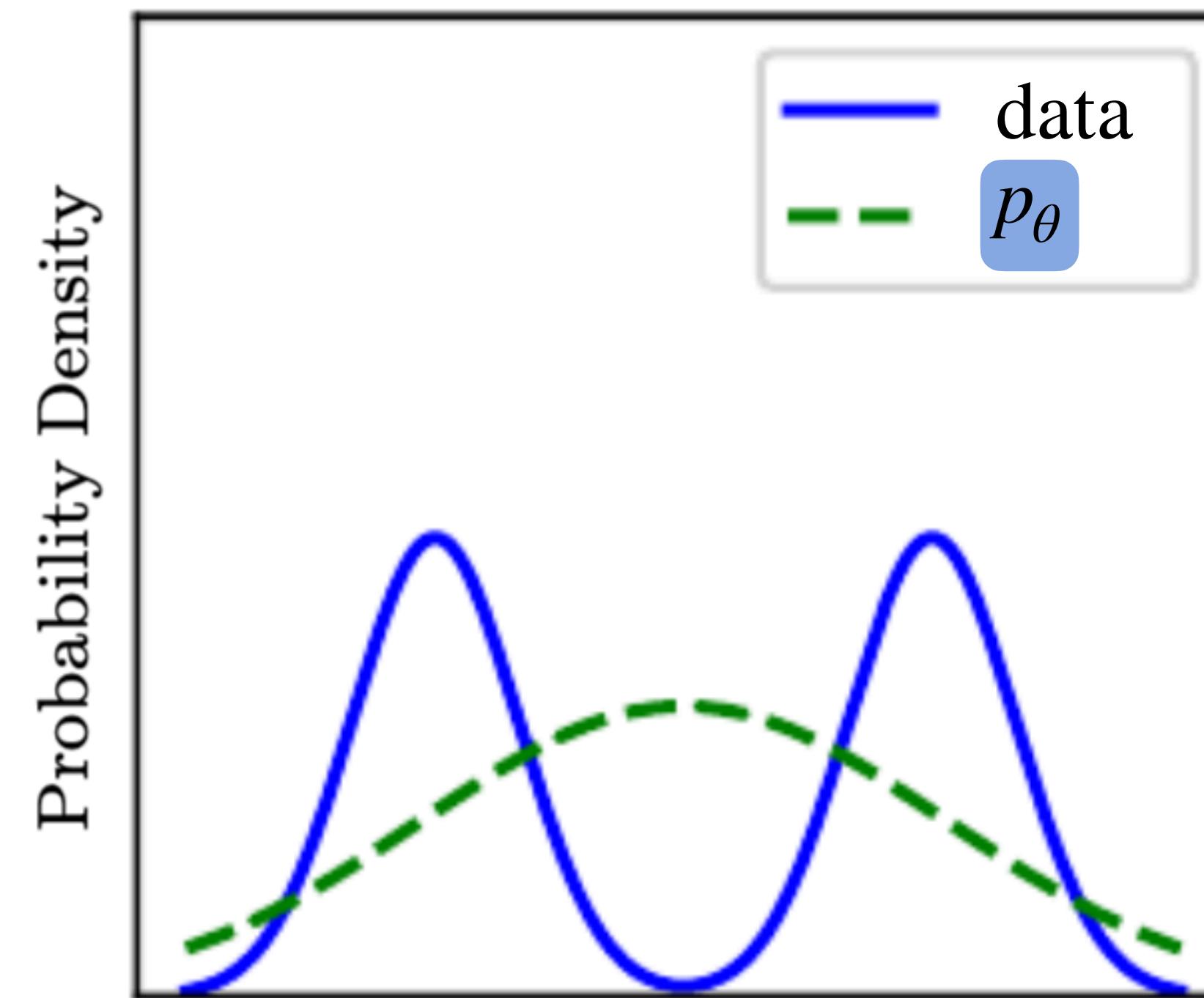


Forward KL or Reverse KL ?

Maximum likelihood estimation

$$\min_{\theta} \text{KL}(\text{data} \parallel p_{\theta})$$

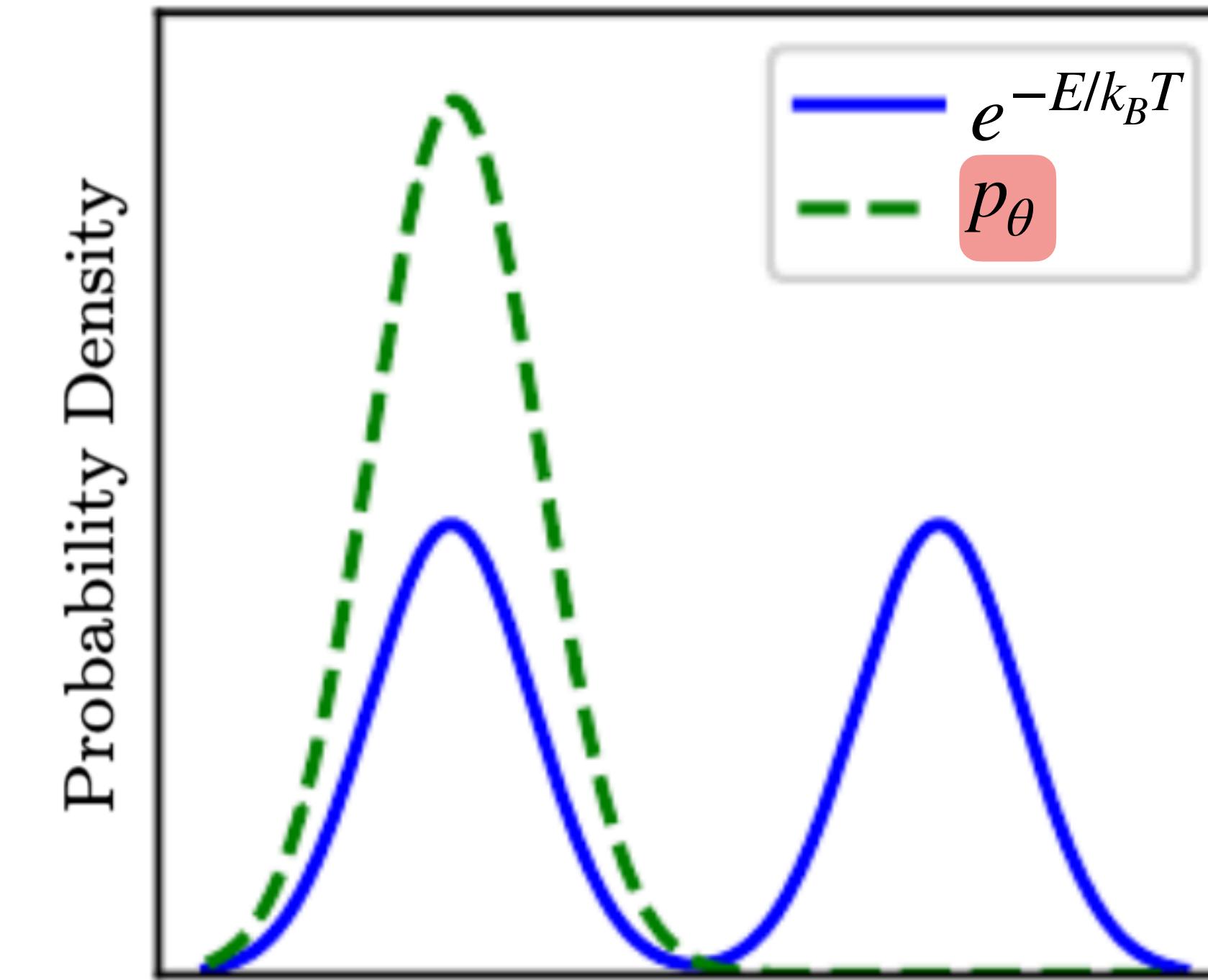
Mode covering

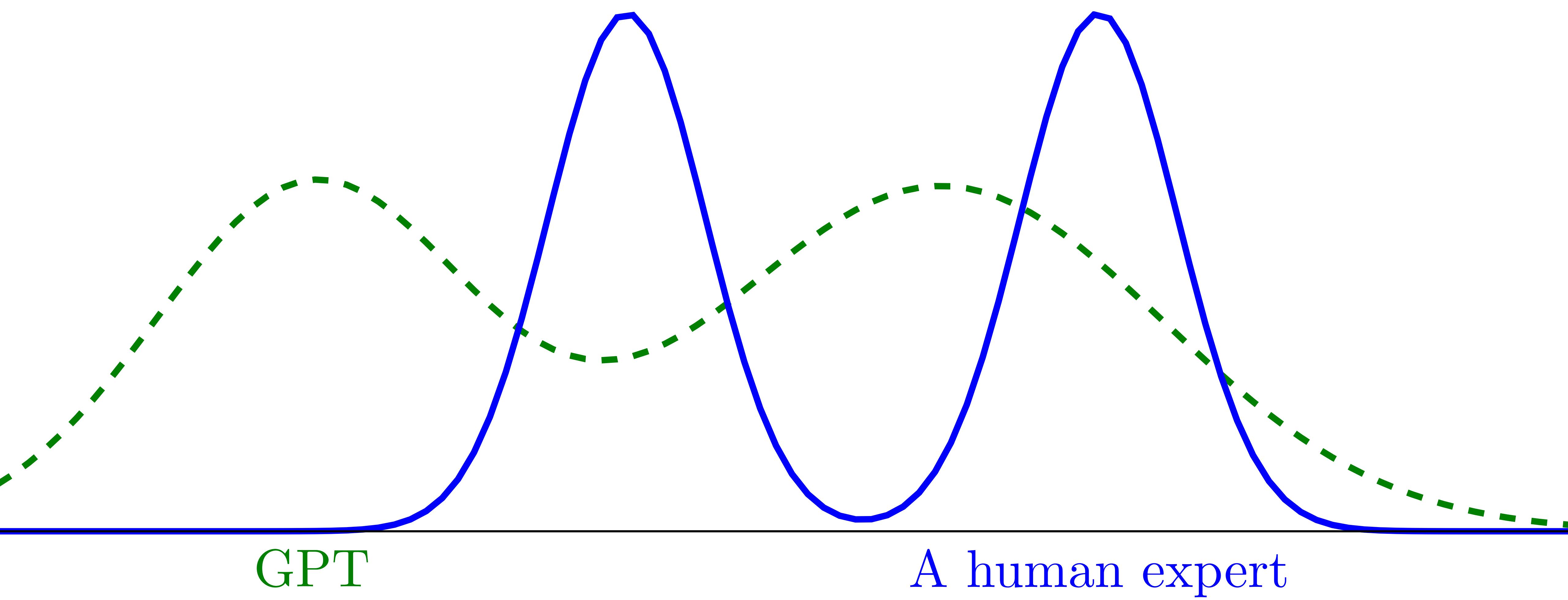


Variational free energy

$$\min_{\theta} \text{KL}(p_{\theta} \parallel e^{-E/k_B T})$$

Mode seeking





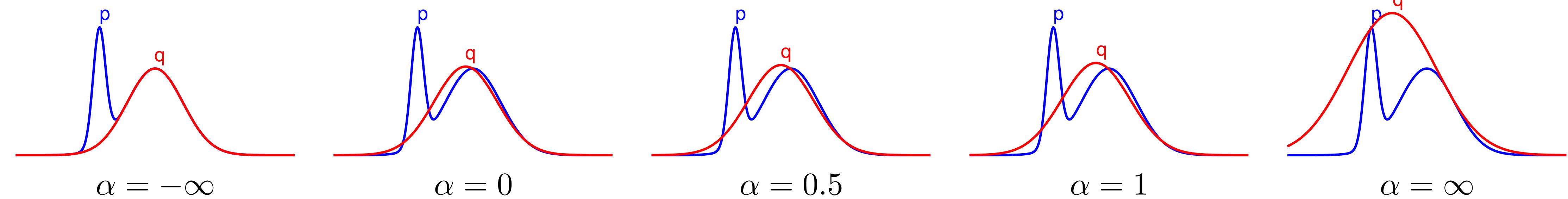
“Jack of all trades, master of none” — 2302.10724

filling the gap vs pushing the boundary of human knowledge

α -divergence

Minka, Microsoft Research Technical Report 2005

$$D_\alpha(p \parallel q) = \frac{\int_x \alpha p(x) + (1 - \alpha)q(x) - p(x)^\alpha q(x)^{1-\alpha} dx}{\alpha(1 - \alpha)}$$



Fisher divergence, defined as

$$F(q, p) = \int_{\mathbb{R}^d} \|\nabla \log q(\theta) - \nabla \log p(\theta)\|^2 q(\theta) d\theta,$$

$$D_{-1}(p \parallel q) = \frac{1}{2} \int_x \frac{(q(x) - p(x))^2}{p(x)} dx$$

$$\lim_{\alpha \rightarrow 0} D_\alpha(p \parallel q) = \text{KL}(q \parallel p)$$

$$D_{\frac{1}{2}}(p \parallel q) = 2 \int_x \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$$

$$\lim_{\alpha \rightarrow 1} D_\alpha(p \parallel q) = \text{KL}(p \parallel q)$$

$$D_2(p \parallel q) = \frac{1}{2} \int_x \frac{(p(x) - q(x))^2}{q(x)} dx$$

Autoregressive models

$$p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)\cdots$$

Language: GPT 2005.14165



Speech: WaveNet 1609.03499

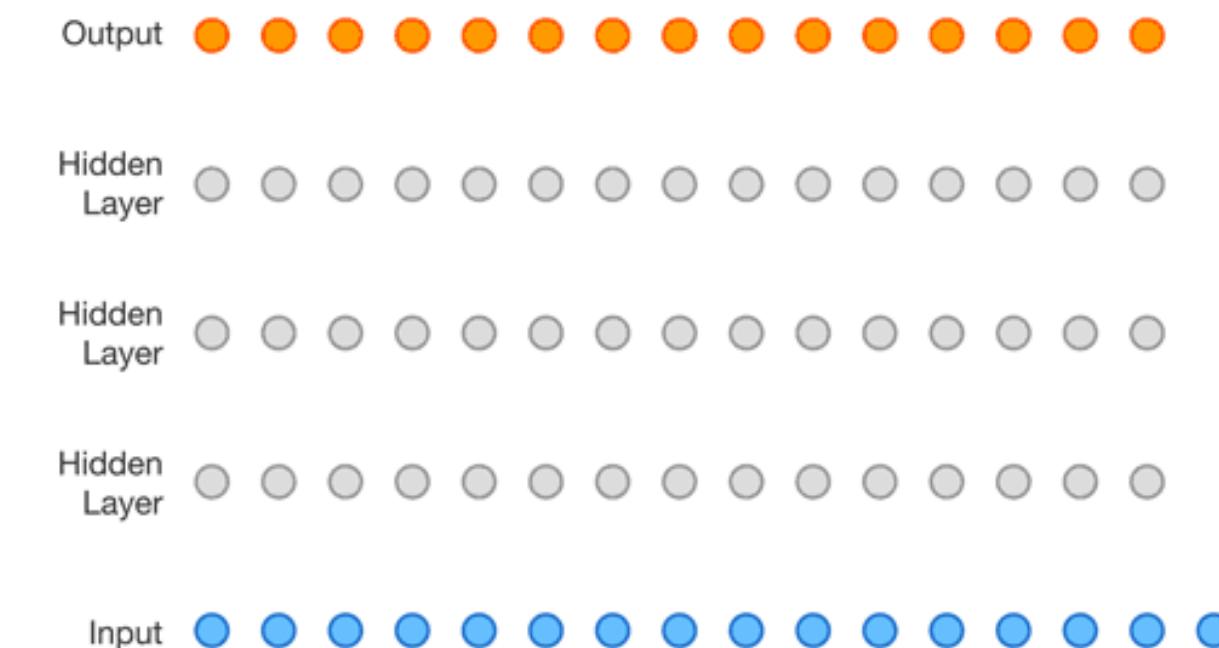
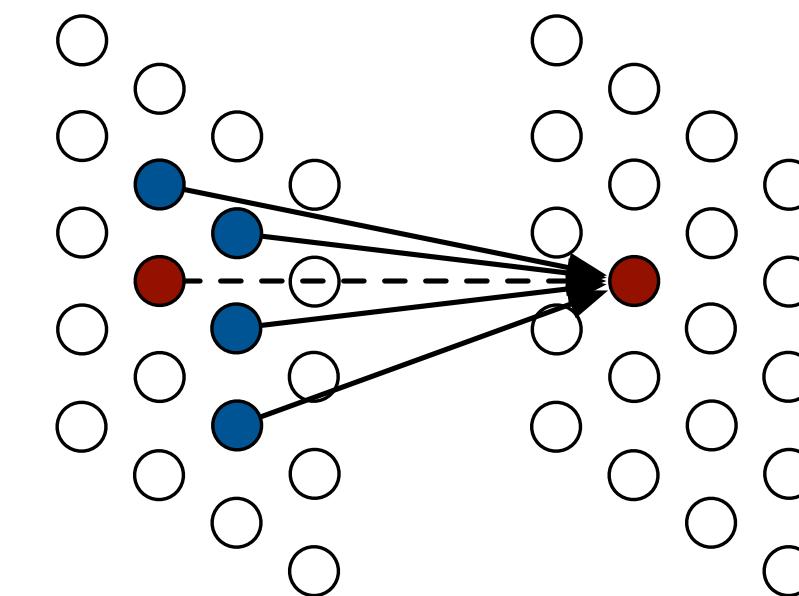
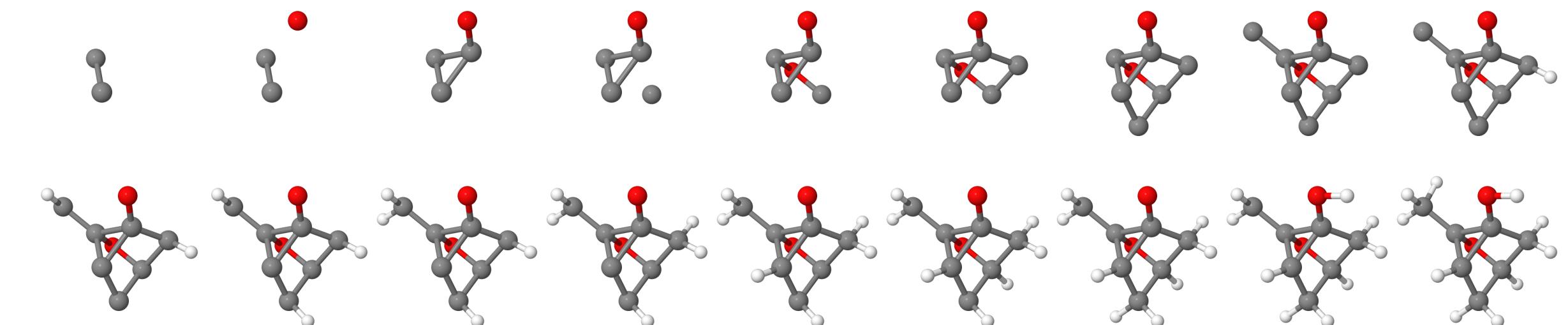


Image: PixelCNN 1601.06759



Molecular graph: 1810.11347



Autoregressive models

$$p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)\cdots$$

Language: GPT 2005.14165



Speech: WaveNet 1609.03499

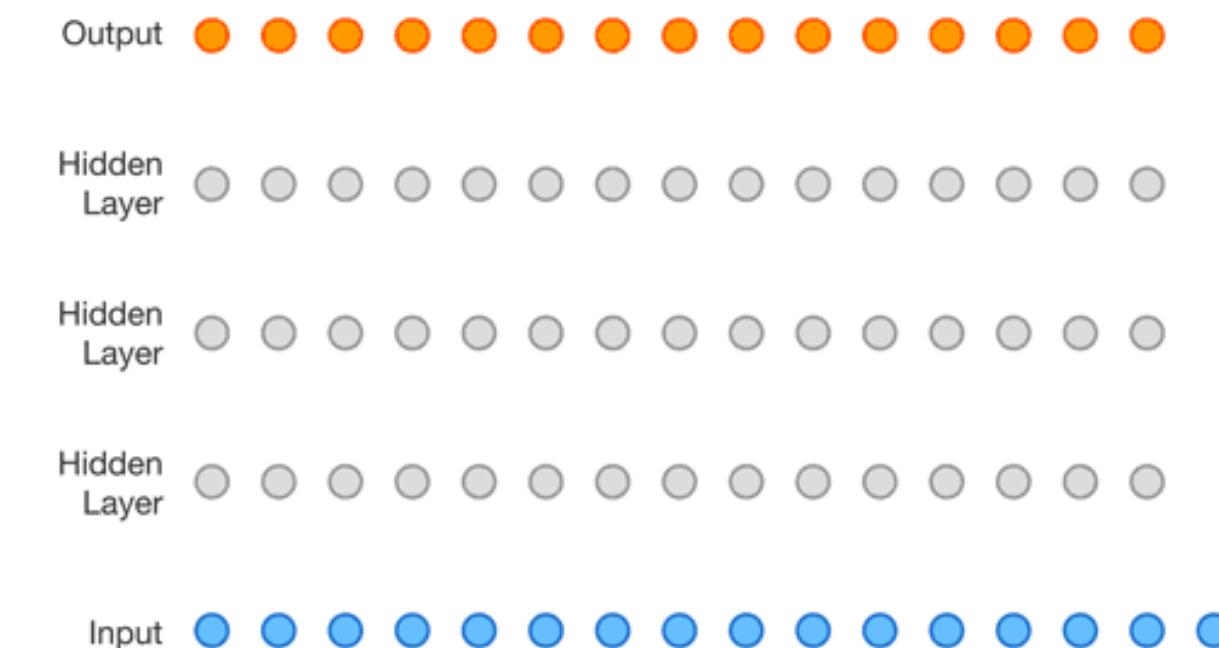
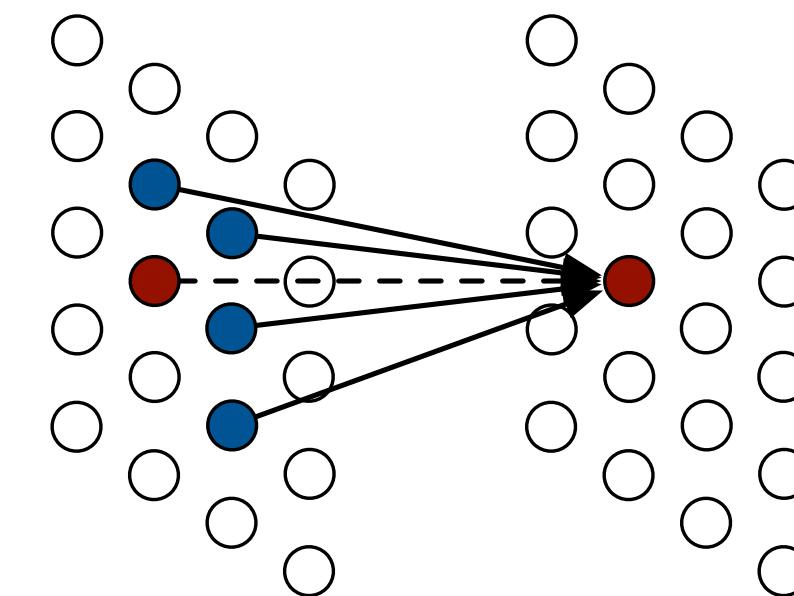
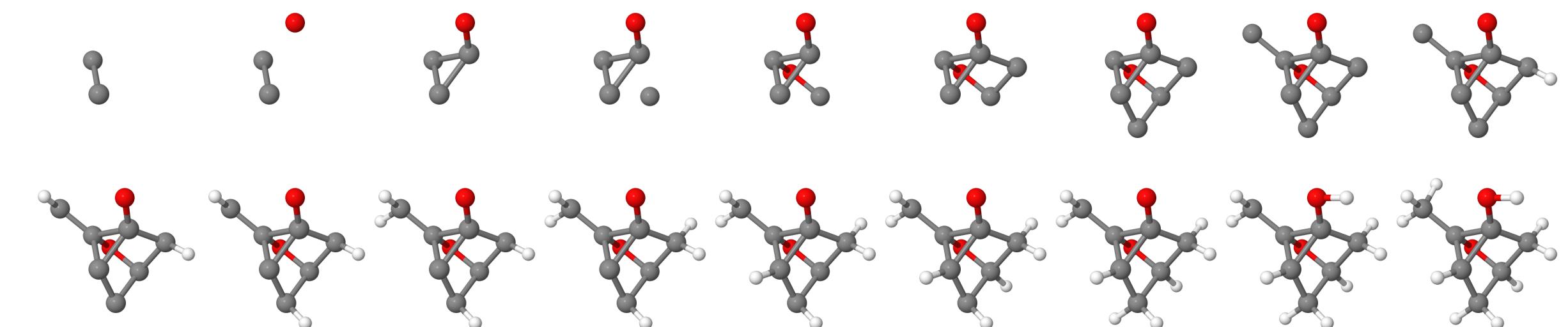


Image: PixelCNN 1601.06759

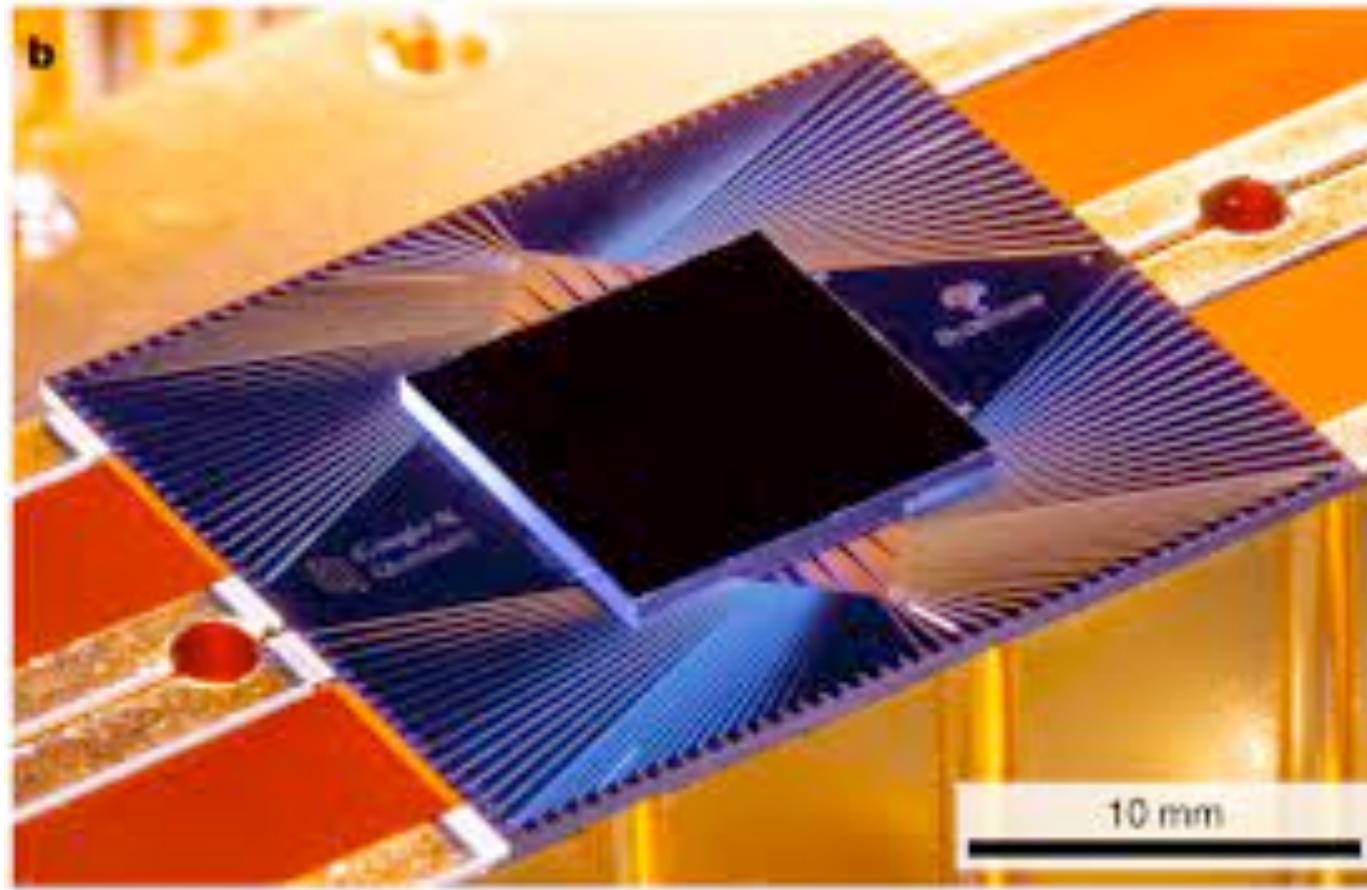


Molecular graph: 1810.11347



Demo: Generative model of Sycamore data

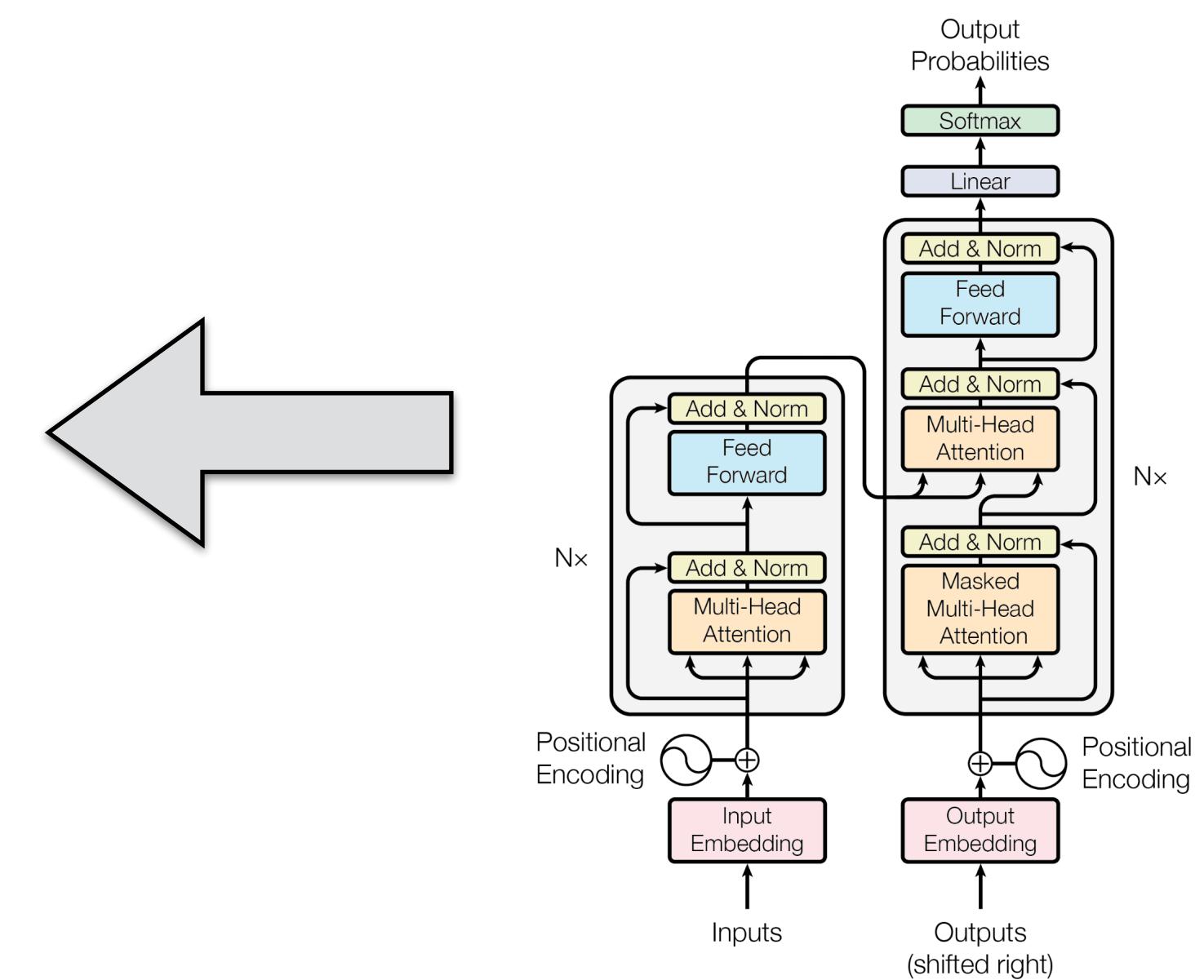
Quantum chip



bitstrings $\sim |\Psi(x)|^2$

```
011110110100  
100001111011  
100110110111  
100110100010  
010100011000  
010001000000  
010101101100  
100001111000  
100101001001  
00100001010
```

Transformer



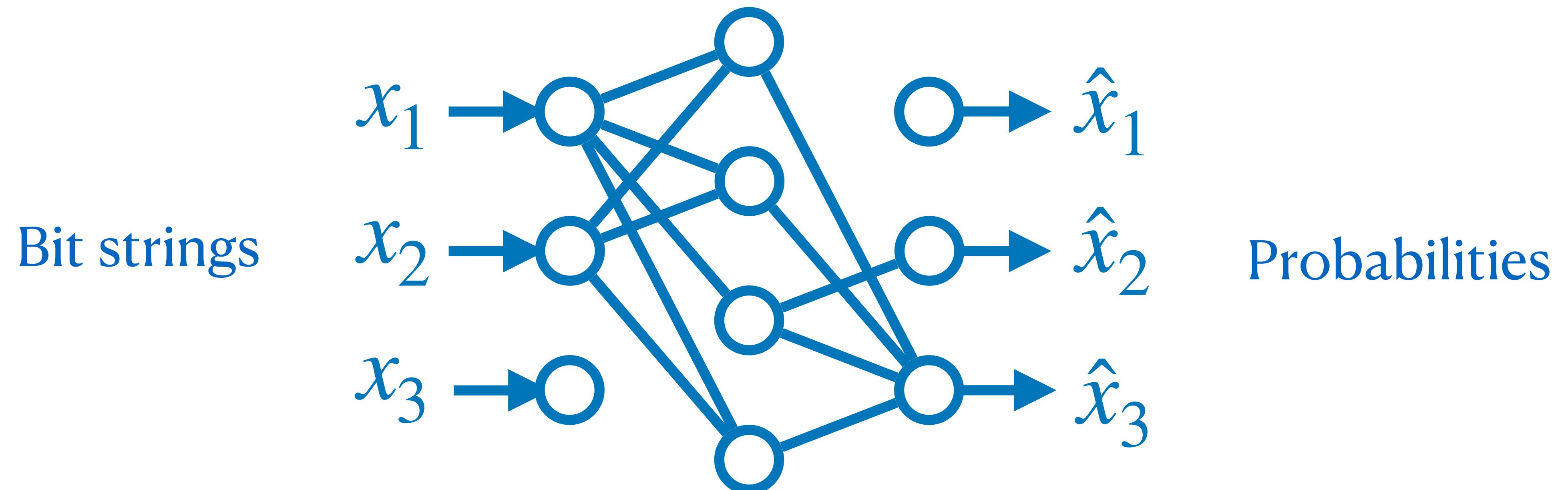
Can we fake the measurement of the sycamore quantum circuit by training a transformer?



https://colab.research.google.com/drive/11WaroqULkudKT3h2i5J6r_EmA4wFKkoZ?usp=sharing

Implementation: autoregressive masks

Masked Autoencoder Germain et al, 1502.03509



$$p(x_1) = \text{Bernoulli}(\hat{x}_1)$$

$$p(x_2 | x_1) = \text{Bernoulli}(\hat{x}_2)$$

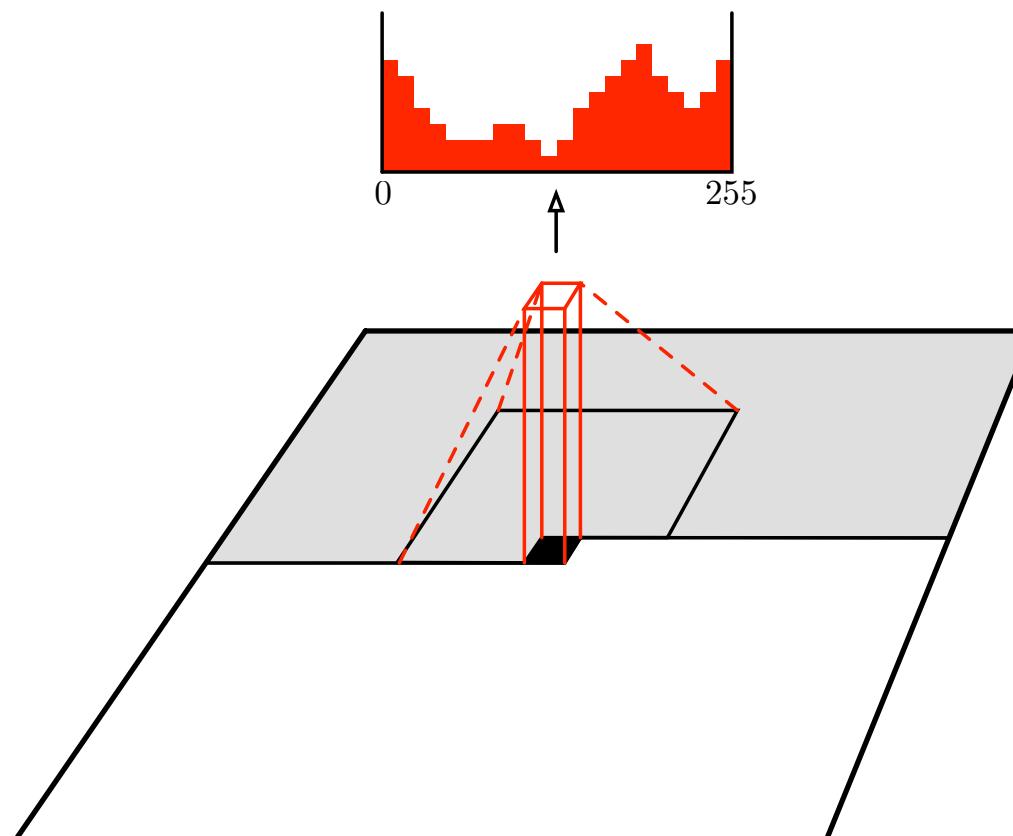
$$p(x_3 | x_1, x_2) = \text{Bernoulli}(\hat{x}_3)$$

Other ways to implement autoregressive models: recurrent networks

Implementation: autoregressive masks

Masked convolutional

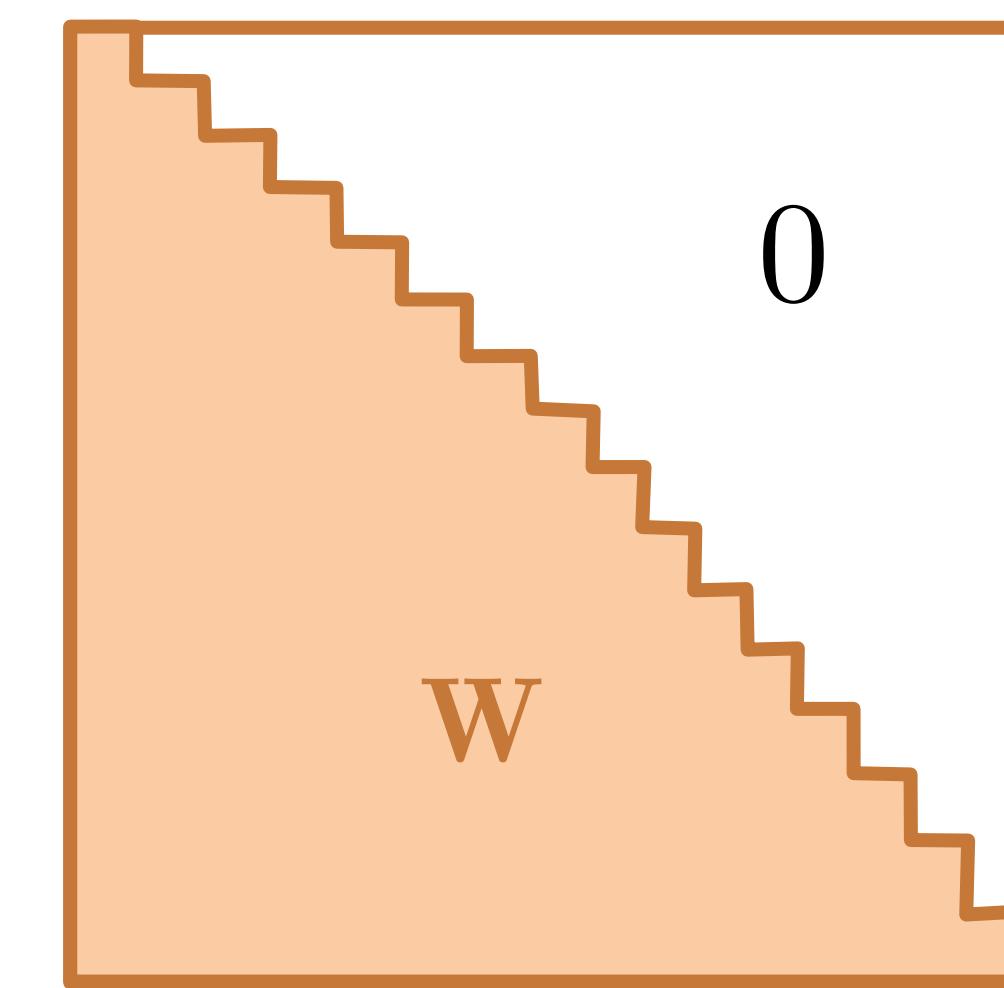
PixelCNN, van den Oord et al, 1601.06759



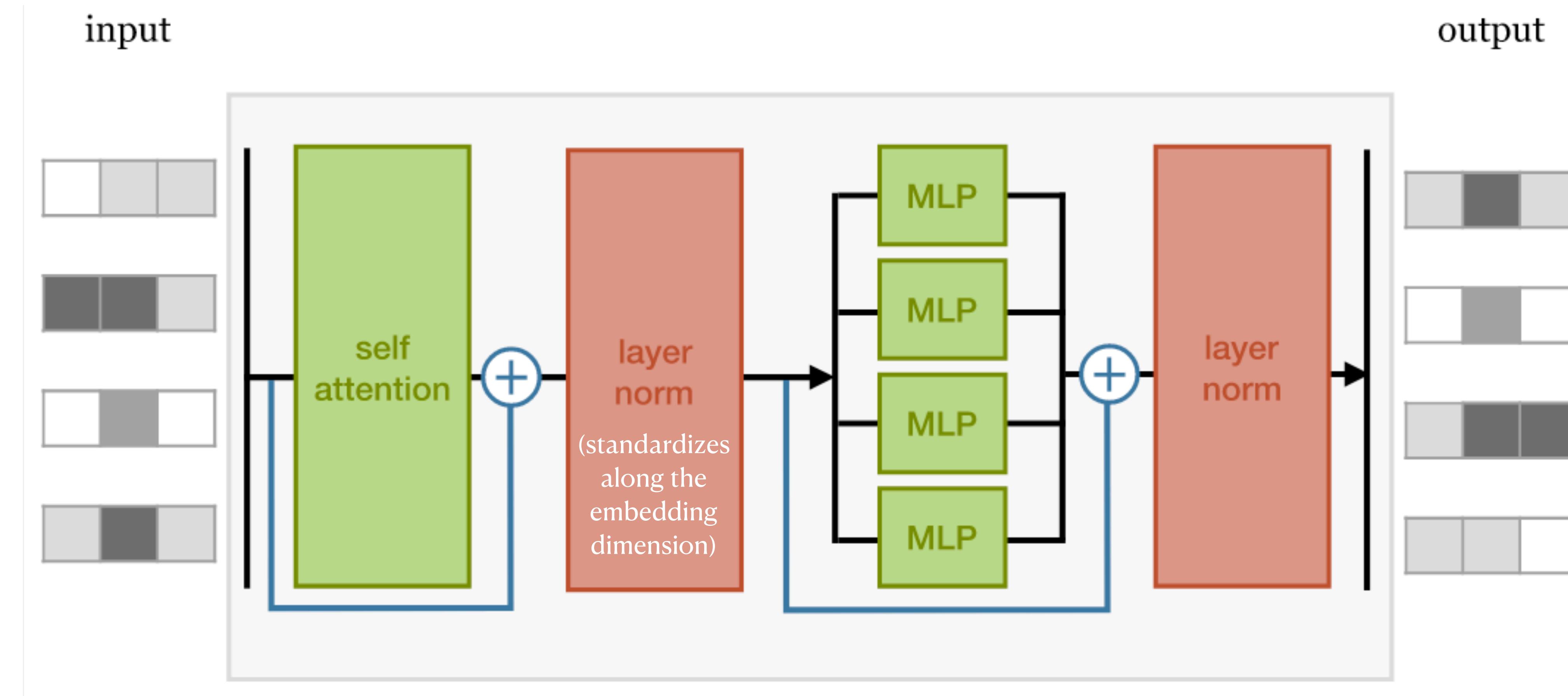
1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

Masked self-attention

Causal transformer, 1706.03762



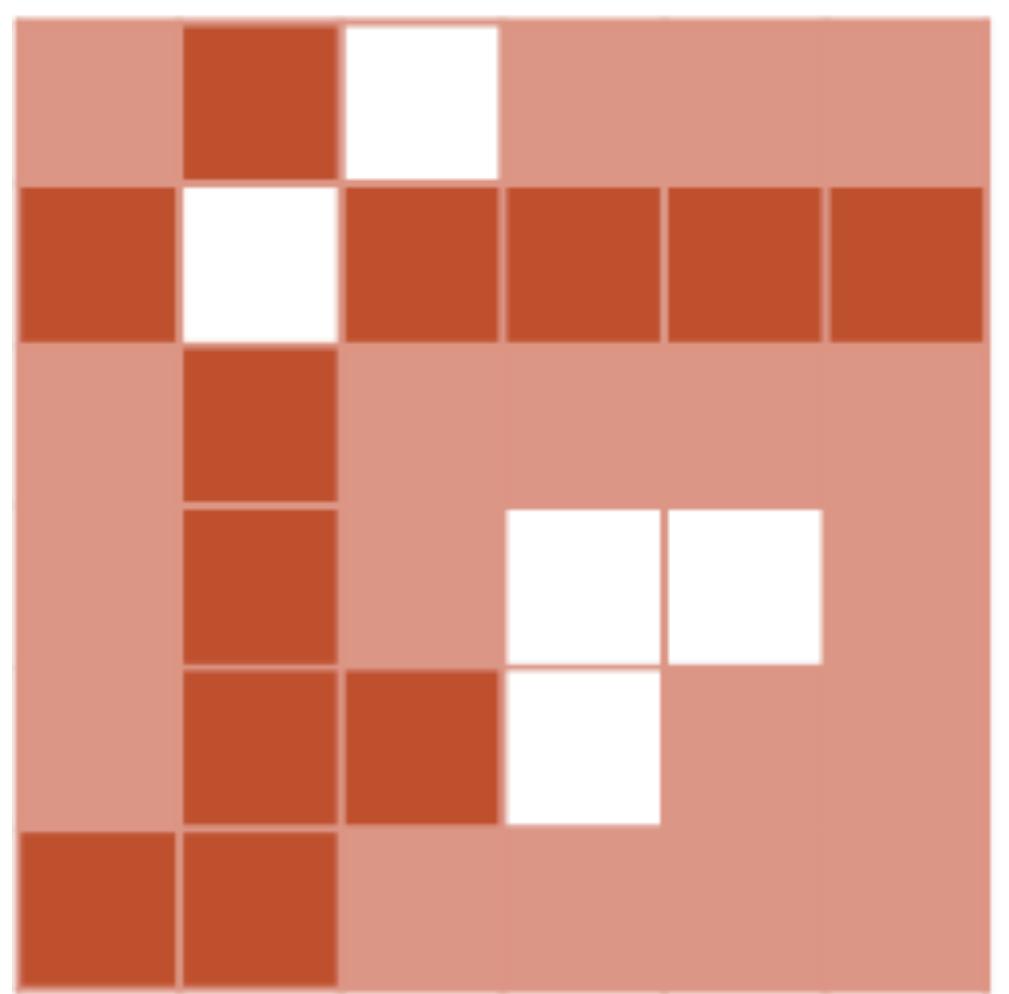
The transformer block



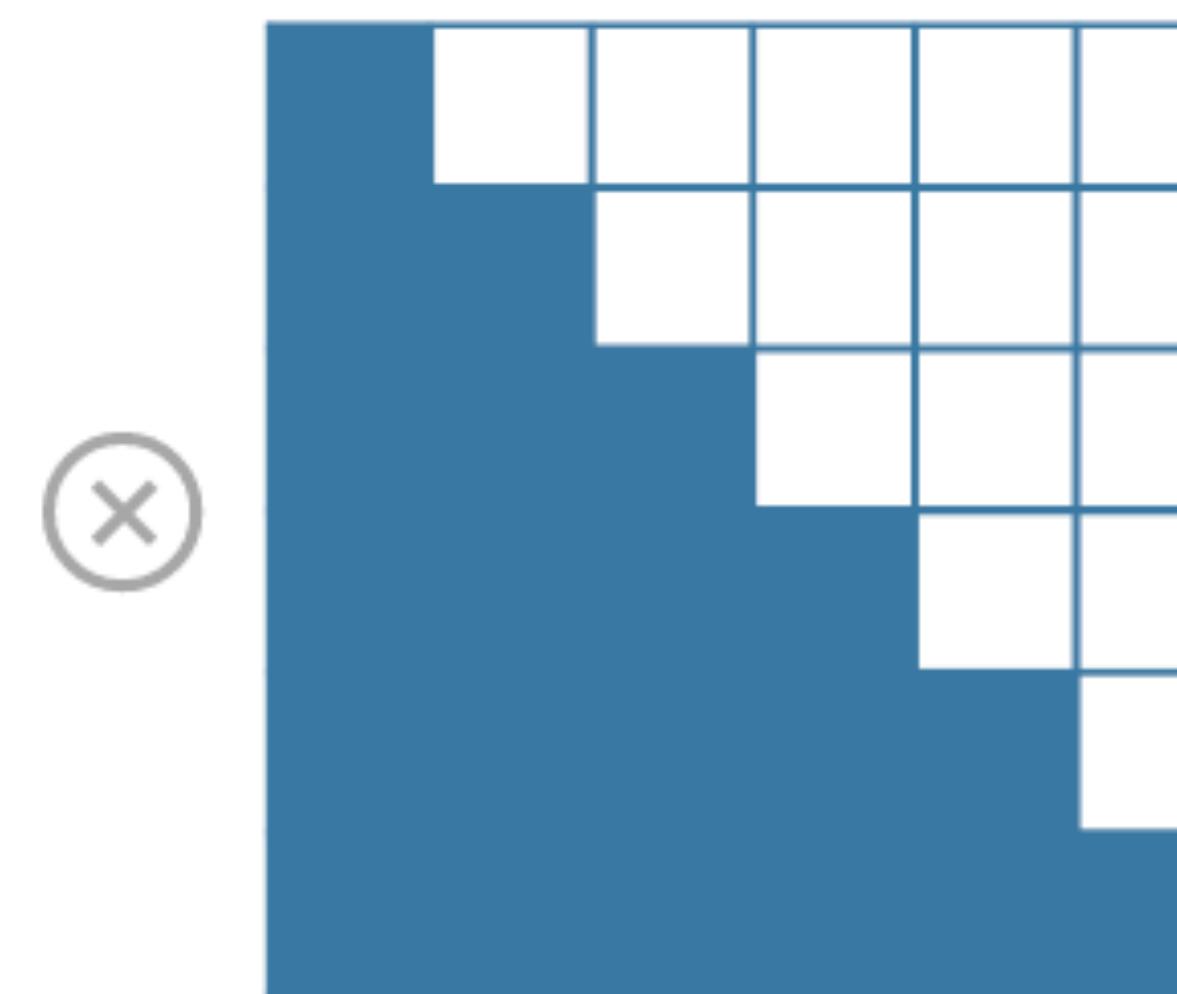
Masked self-attention

$$y_i = \sum_j \alpha(x_i, x_j) x_j$$

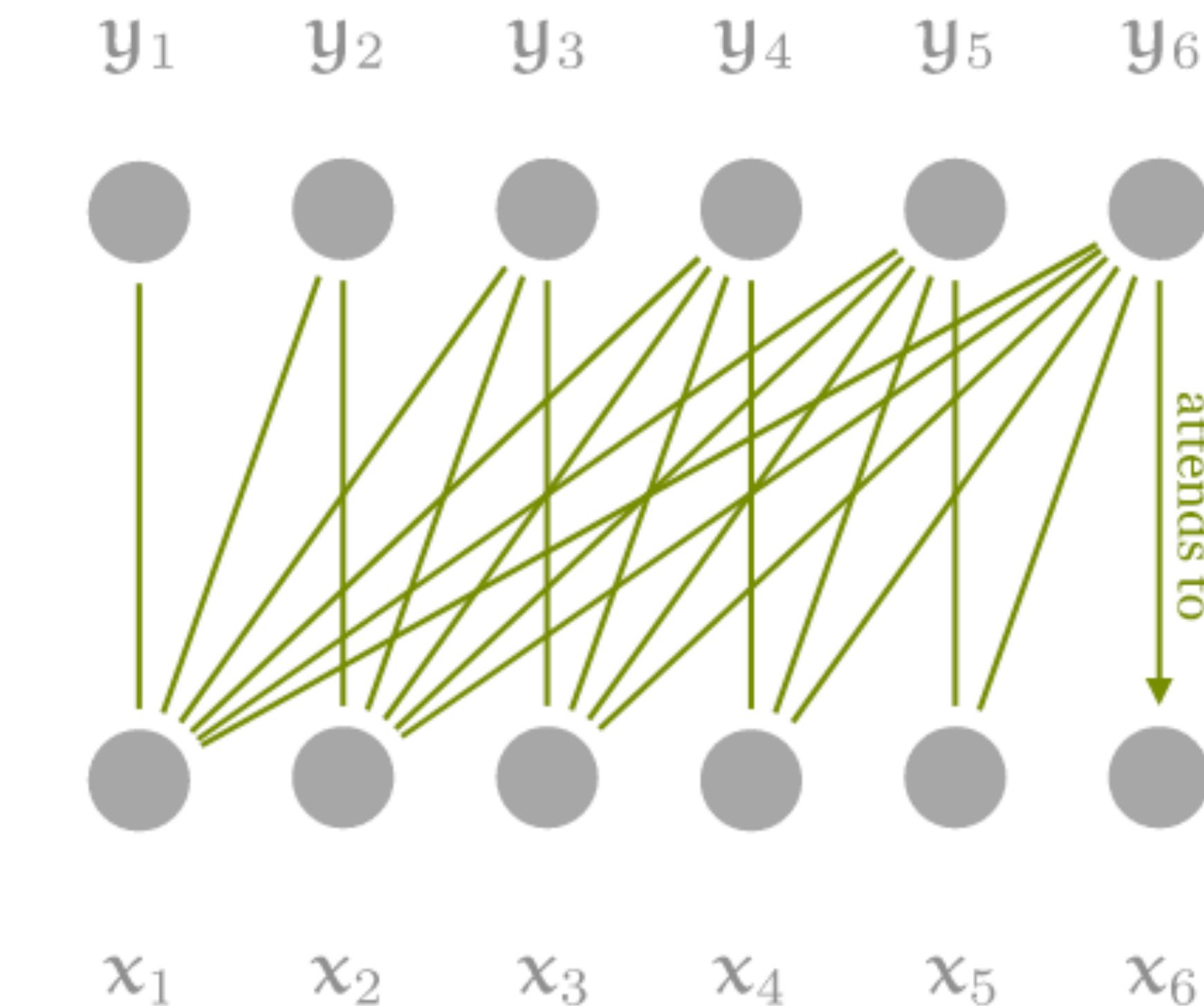
↗ attention weight



raw attention weights



mask

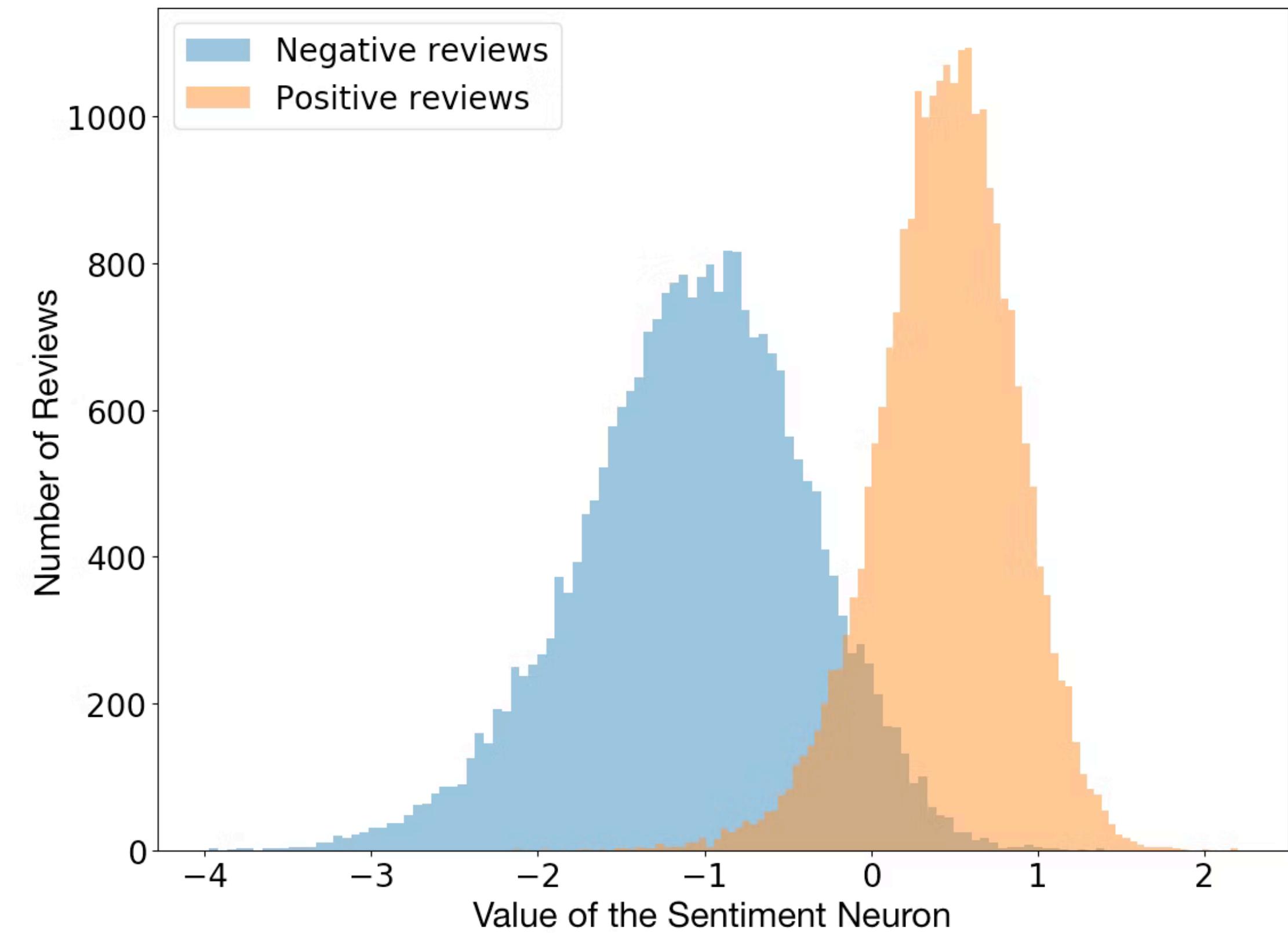


Learning to Generate Reviews and Discovering Sentiment

Alec Radford¹ Rafal Jozefowicz¹ Ilya Sutskever¹

We explore the properties of byte-level recurrent language models. When given sufficient amounts of capacity, training data, and compute time, the representations learned by these models include disentangled features corresponding to high-level concepts. Specifically, we find a single unit which performs sentiment analysis. These representations, learned in an unsupervised manner, achieve state of the art on the binary subset of the Stanford Sentiment Treebank. They are also very data efficient. When using only a handful of labeled examples, our approach matches the performance of strong baselines trained on full datasets. We also demonstrate the sentiment unit has a direct influence on the generative process of the model. Simply fixing its value to be positive or negative generates samples with the corresponding positive or negative sentiment.

“Sentiment neuron”



Generative Pretraining from Pixels

Mark Chen¹ Alec Radford¹ Rewon Child¹ Jeff Wu¹ Heewoo Jun¹ Prafulla Dhariwal¹ David Luan¹
Ilya Sutskever¹

Inspired by progress in unsupervised representation learning for natural language, we examine whether similar models can learn useful representations for images. We train a sequence Transformer to auto-regressively predict pixels, without incorporating knowledge of the 2D input structure. Despite training on low-resolution ImageNet without labels, we find that a GPT-2 scale model learns strong image representations as measured by linear probing, fine-tuning, and low-data classification. On CIFAR-10, we achieve 96.3% accuracy with a linear probe, outperforming a supervised Wide ResNet, and 99.0% accuracy with full fine-tuning, matching the top supervised pre-trained models. An even larger model trained on a mixture of ImageNet and web images is competitive with self-supervised benchmarks on ImageNet, achieving 72.0% top-1 accuracy on a linear probe of our features.

Representation learned by image GPT

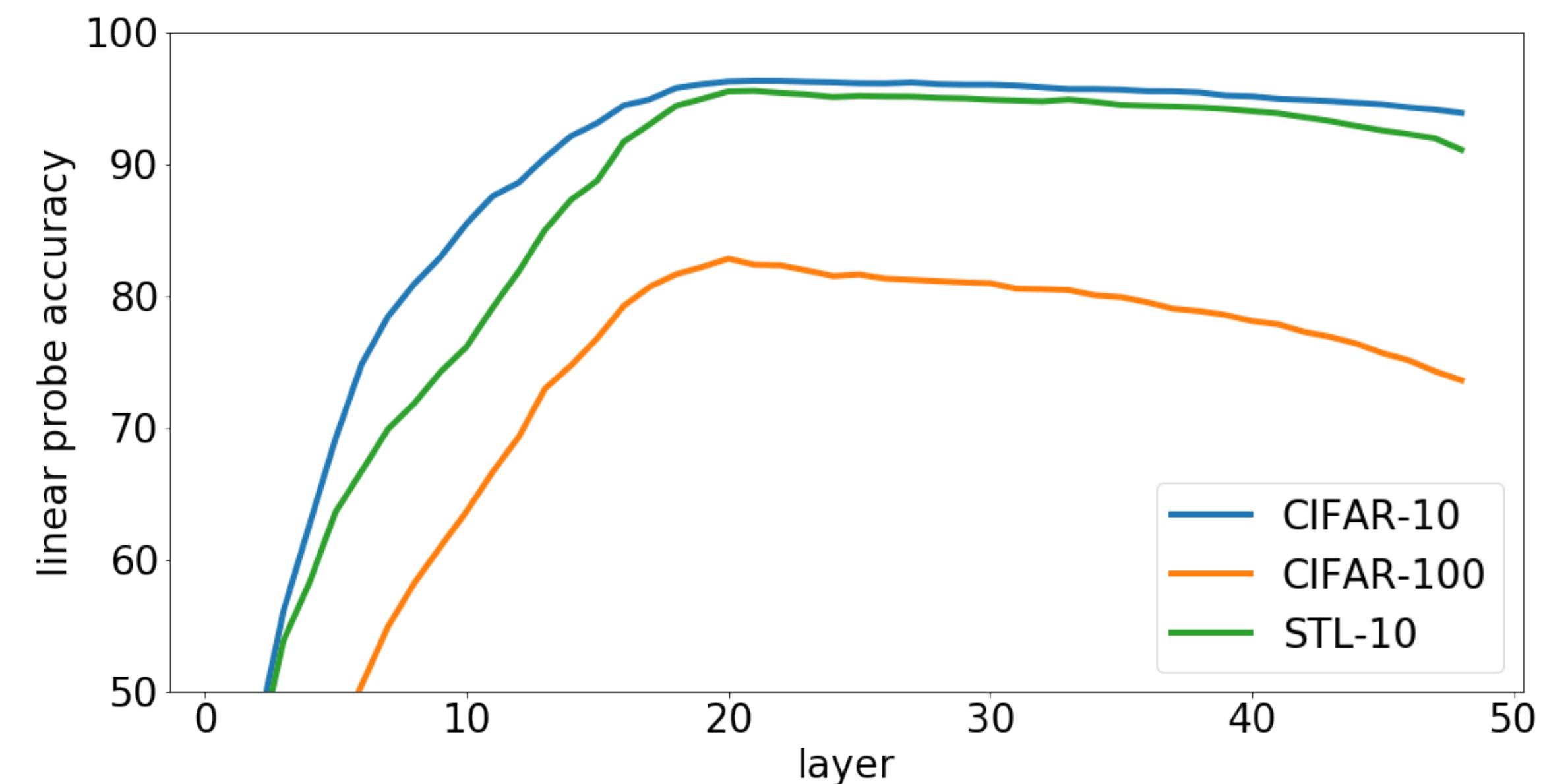
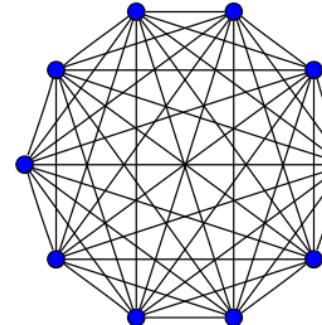
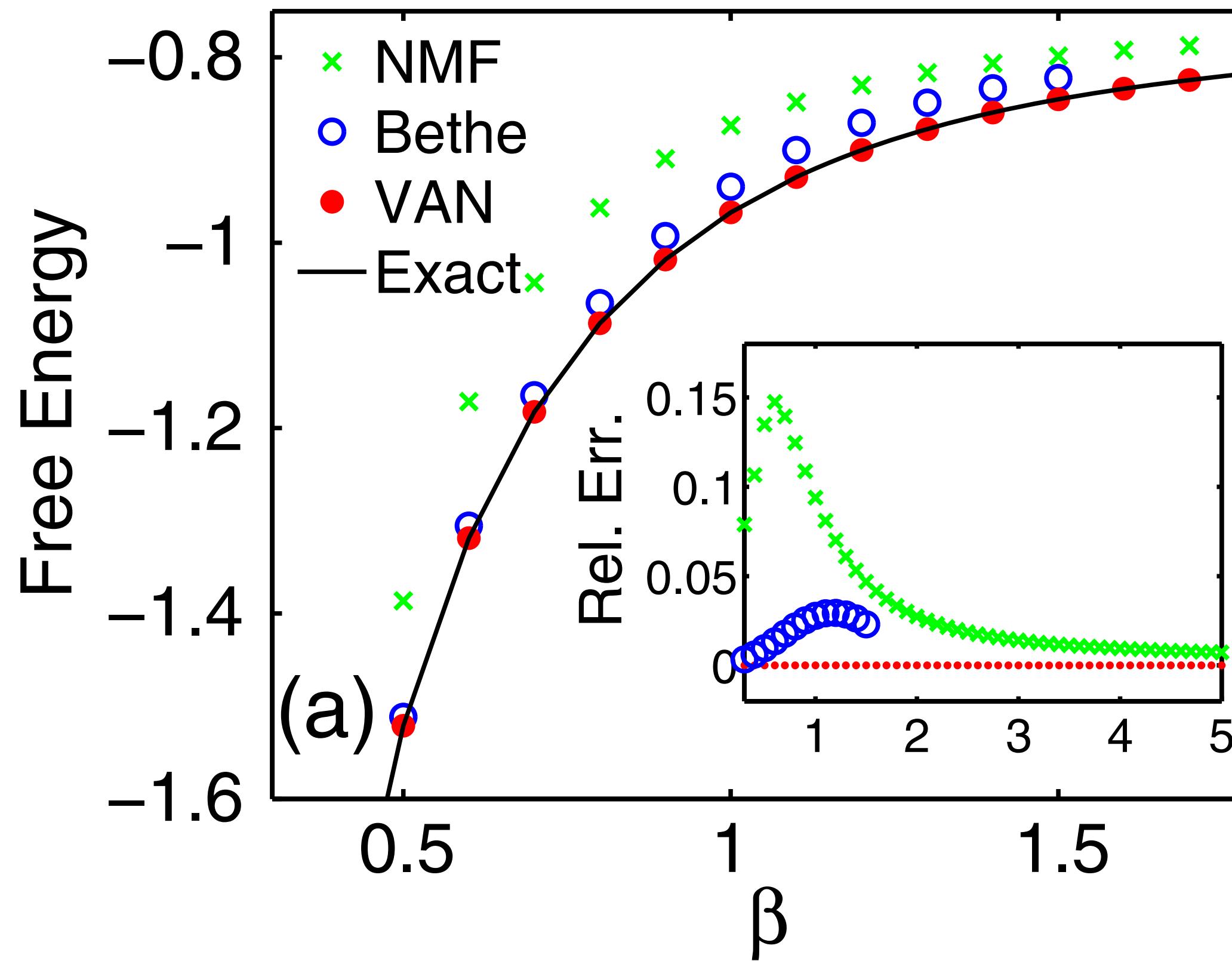


Figure 2. Representation quality depends on the layer from which we extract features. In contrast with supervised models, the best representations for these generative models lie in the middle of the network. We plot this unimodal dependence on depth by showing linear probes for iGPT-L on CIFAR-10, CIFAR-100, and STL-10.

Variational autoregressive network for statistical mechanics



Sherrington-Kirkpatrick spin glass



Naive mean-field
factorized probability

$$p(X) = \prod_i p(x_i)$$

Bethe approximation
pairwise interaction

$$p(X) = \prod_i p(x_i) \prod_{(i,j) \in E} \frac{p(x_i, x_j)}{p(x_i)p(x_j)}$$

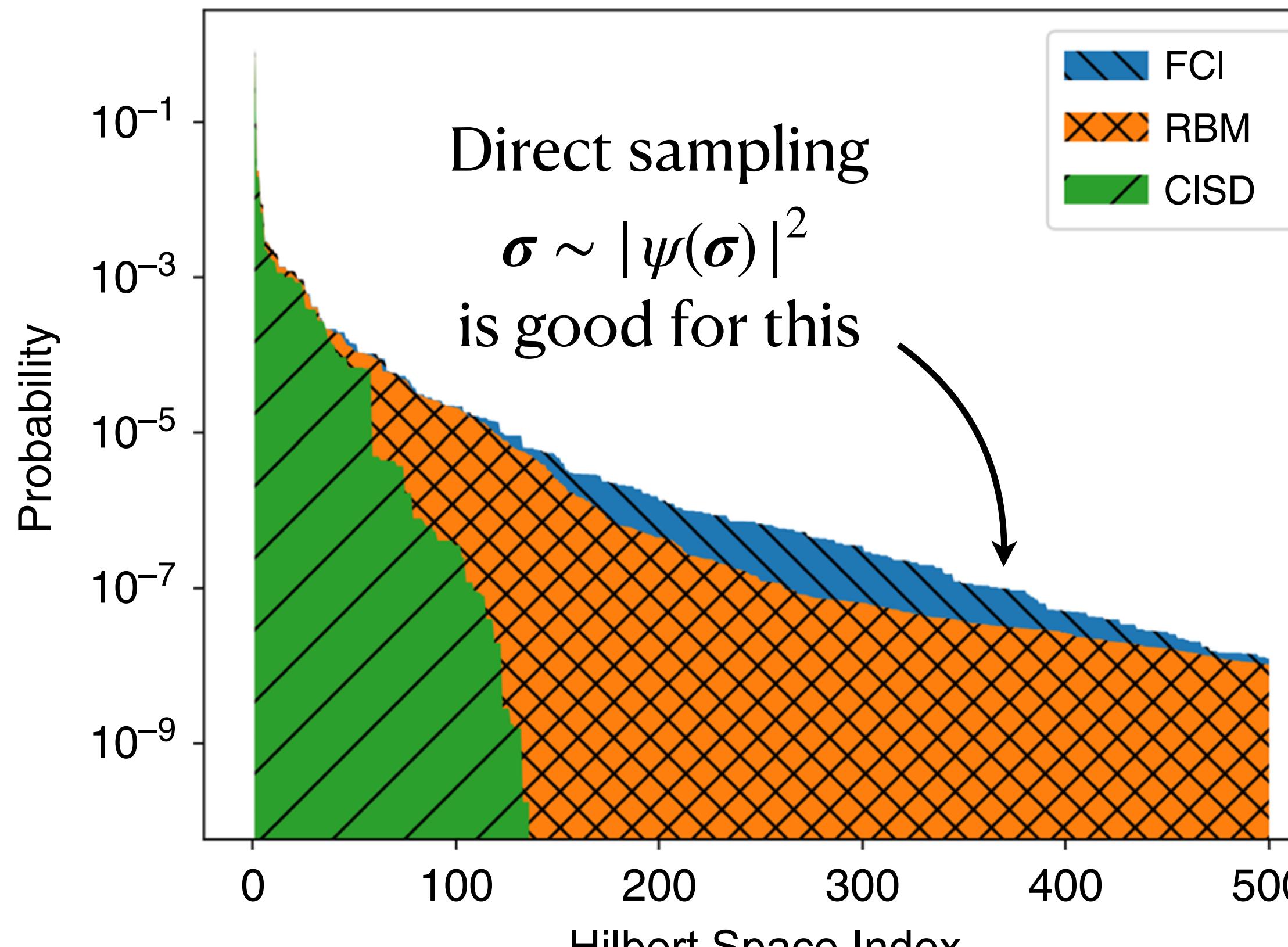
Variational autoregressive
network

$$p(X) = \prod_i p(x_i | x_{<i})$$

Wu, LW, Zhang, PRL '19
github.com/wdphy16/stat-mech-van

Variational autoregressive quantum states

$$\psi(\sigma) = \psi(\sigma_1)\psi(\sigma_2 | \sigma_1)\psi(\sigma_3 | \sigma_1, \sigma_2)\cdots$$



Objective function: ground state energy

McMillan 1965, Carleo & Troyer Science 2017

$$\frac{\langle \psi | \hat{H} | \psi \rangle}{\langle \psi | \psi \rangle} = \mathbb{E}_{\sigma \sim |\psi(\sigma)|^2} \left[\frac{\hat{H}\psi(\sigma)}{\psi(\sigma)} \right]$$

Sharir, Levine, Wies, Carleo, Shashua, PRL '20

Hibat-Allah, Ganahl, Hayward, Melko, Carrasquilla, PRResearch '20

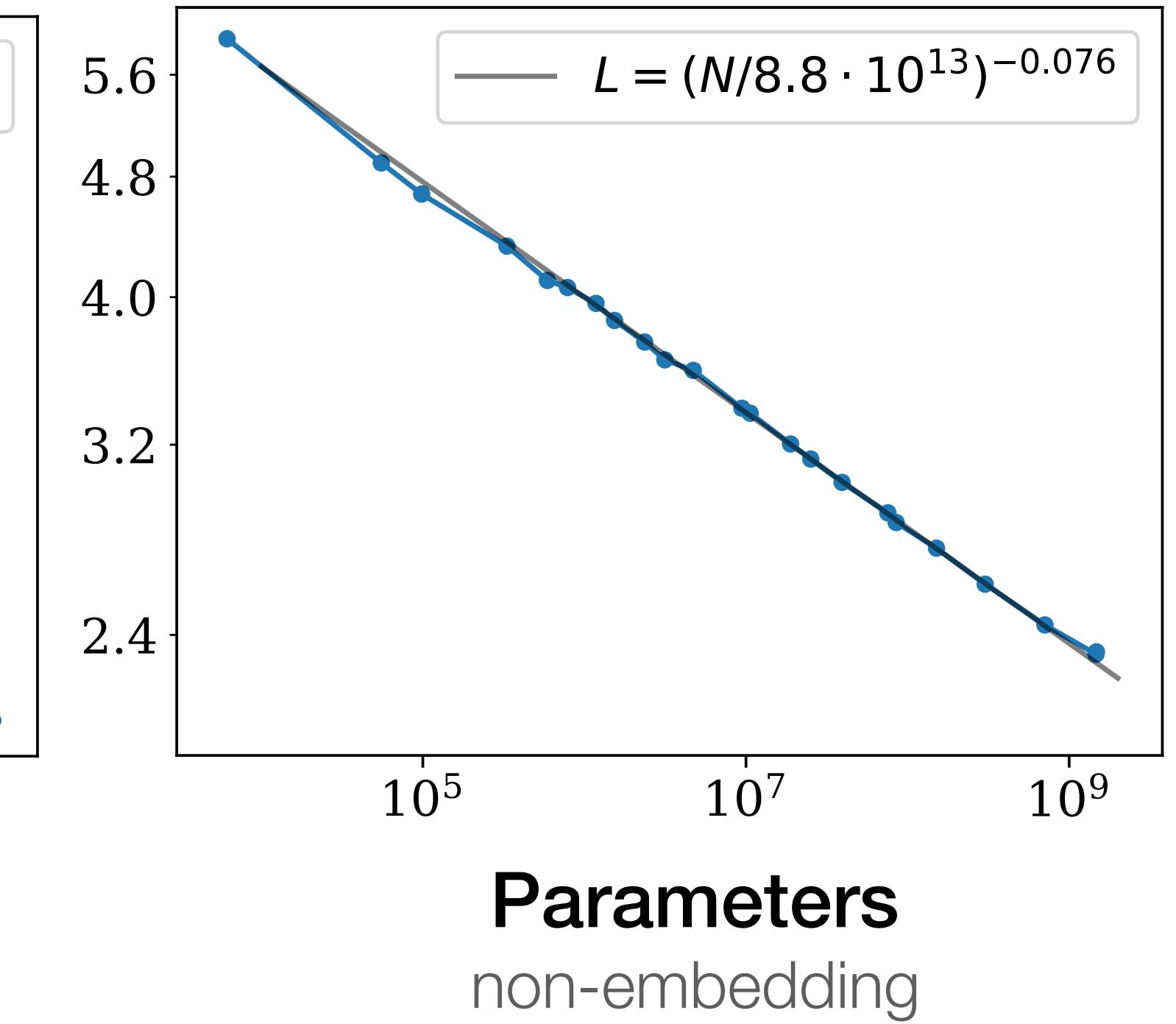
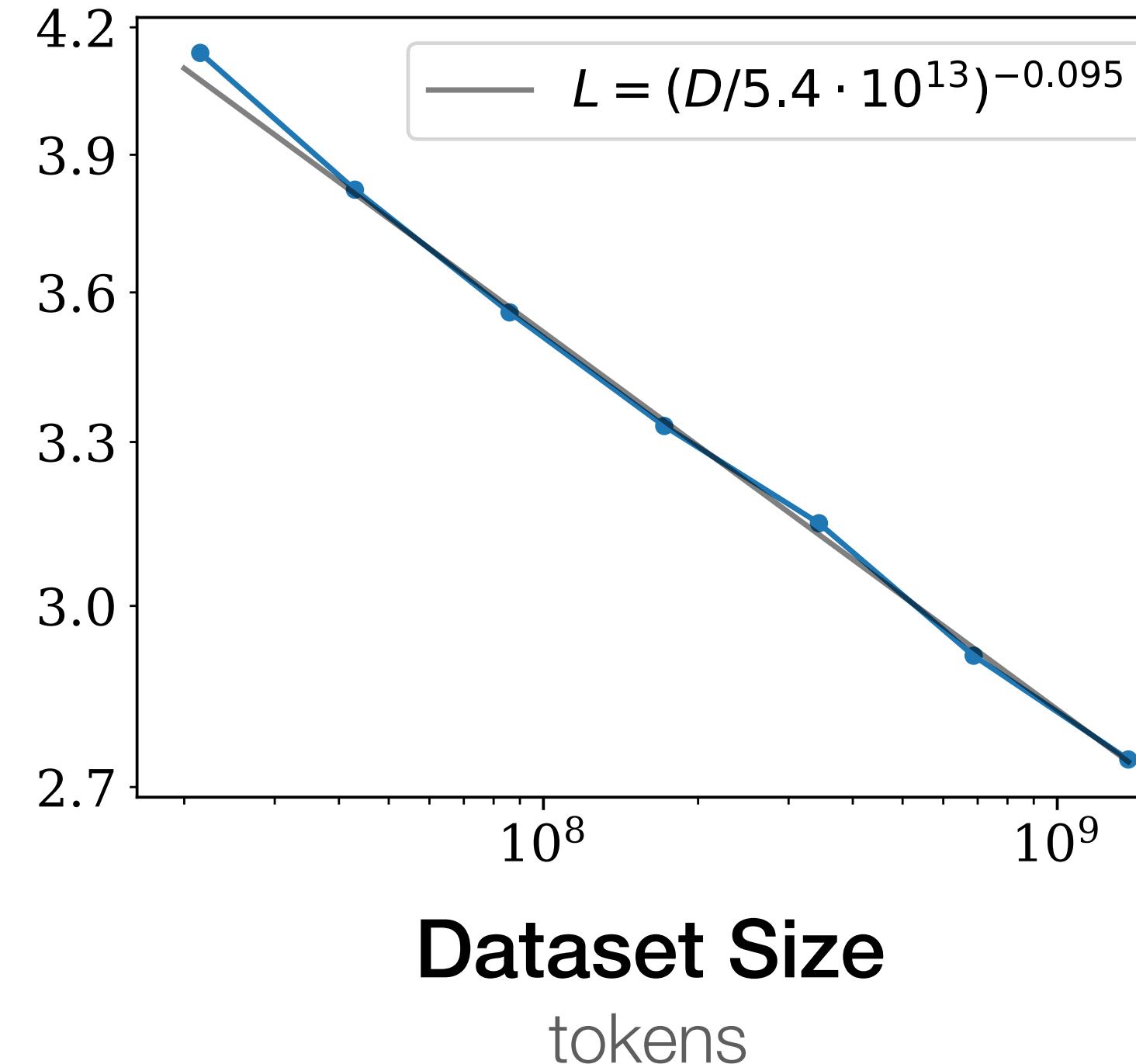
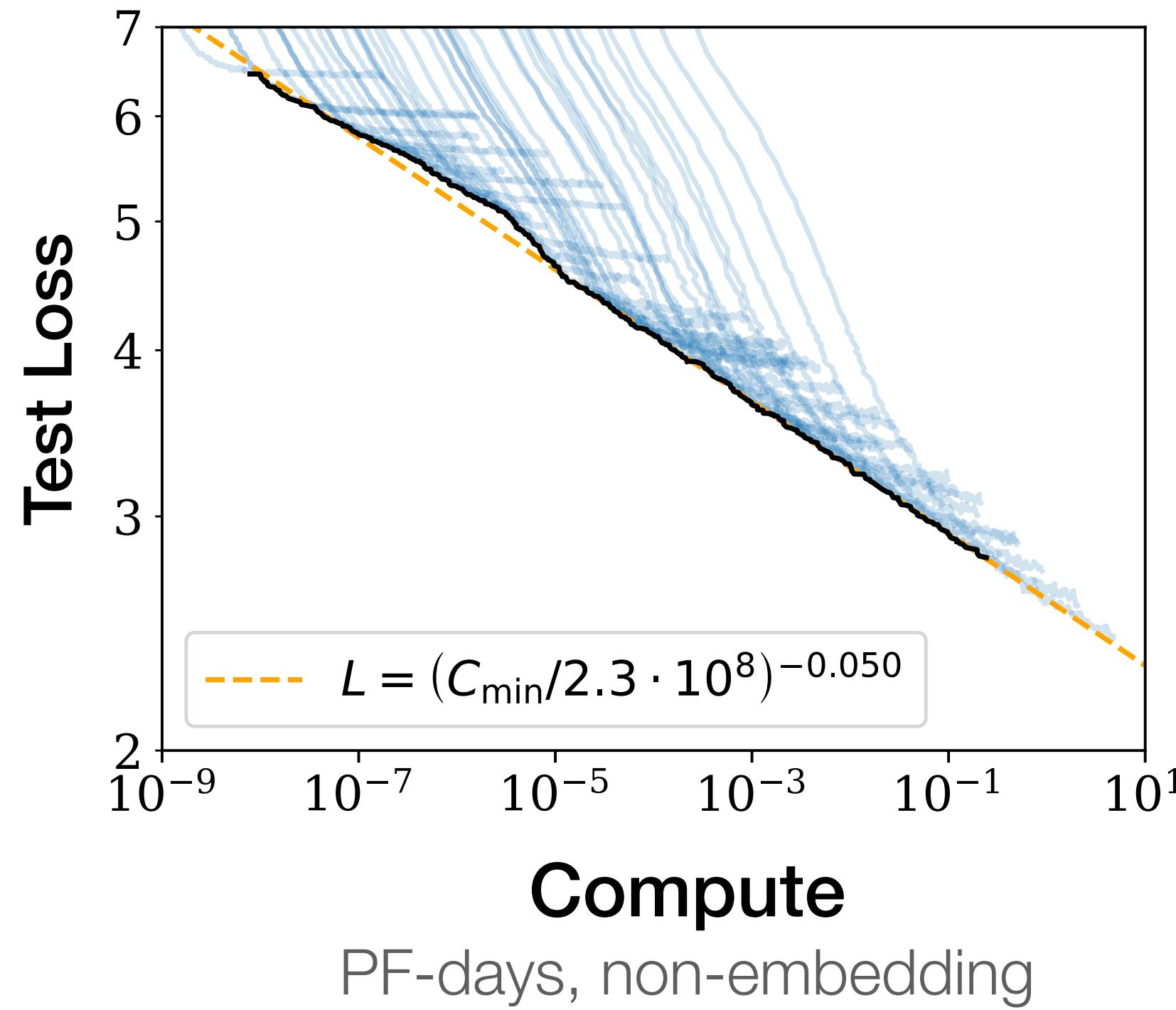
Barrett et al, Nat. Mach. Intell. '22

Zhao et al, MLST. '23

Shang et al, 2307.09343

Scaling law

Kaplan et al, 2001.08361



“It would also be exciting to find a theoretical framework from which the scaling relations can be derived: a ‘statistical mechanics’ underlying the ‘thermodynamics’ we have observed.”

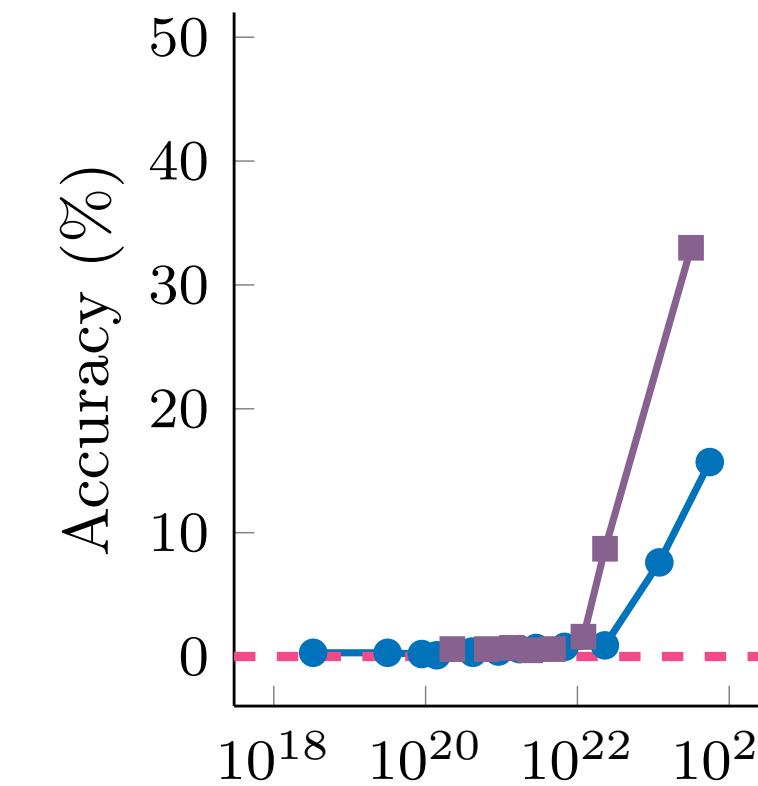
Emergent abilities: more is different

Wei et al, 2206.07682

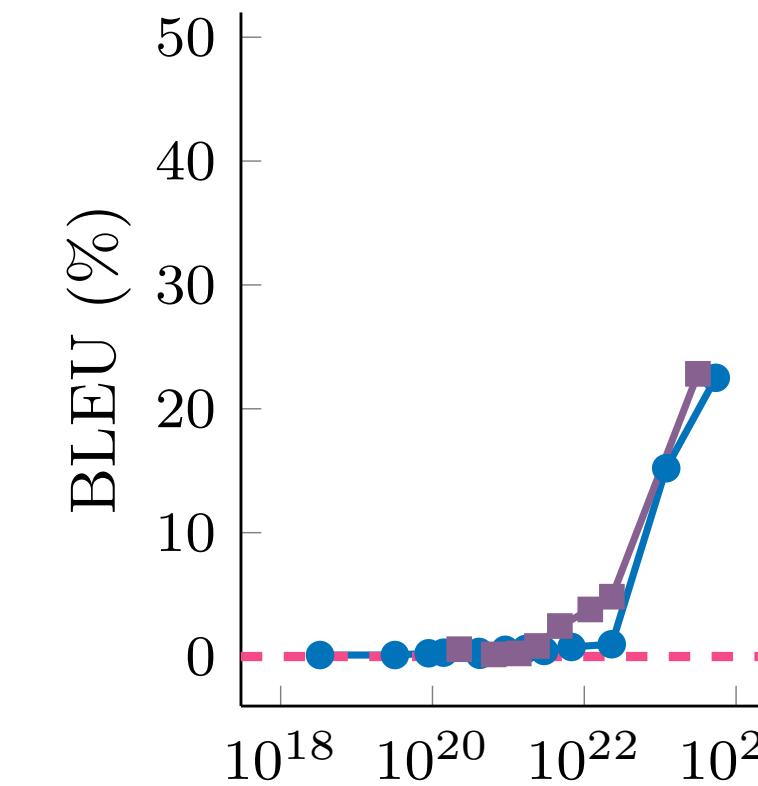
[https://www.jasonwei.net/
blog/emergence](https://www.jasonwei.net/blog/emergence)

● LaMDA ■ GPT-3 ◆ Gopher ▲ Chinchilla ▽ PaLM - - - Random

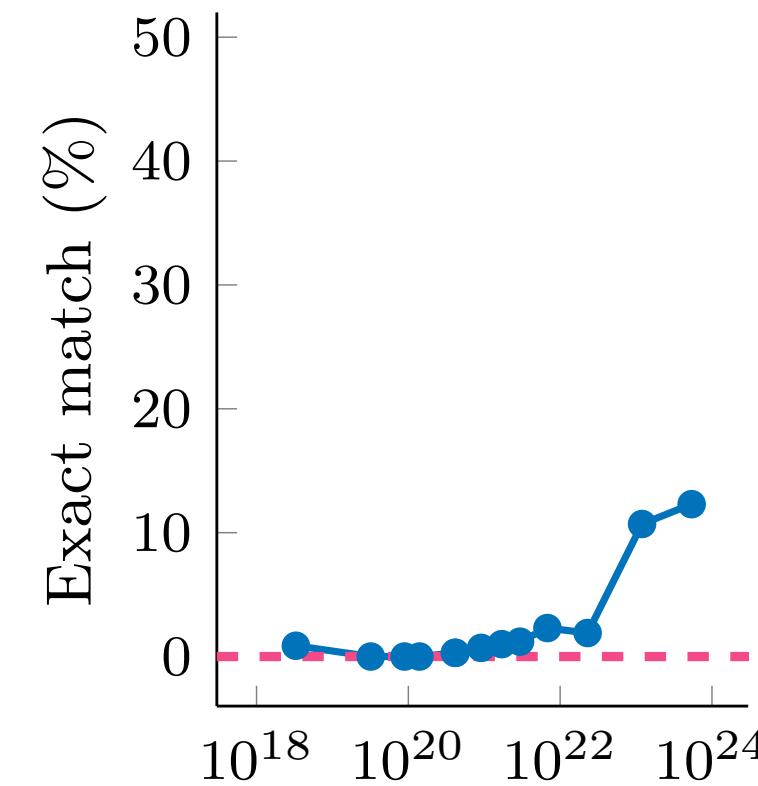
(A) Mod. arithmetic



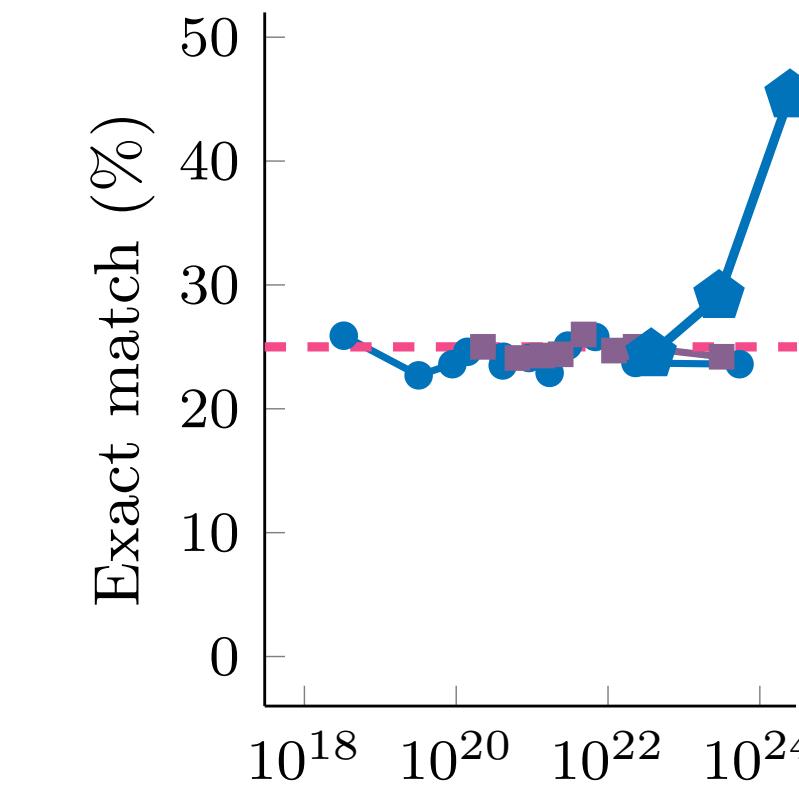
(B) IPA transliterate



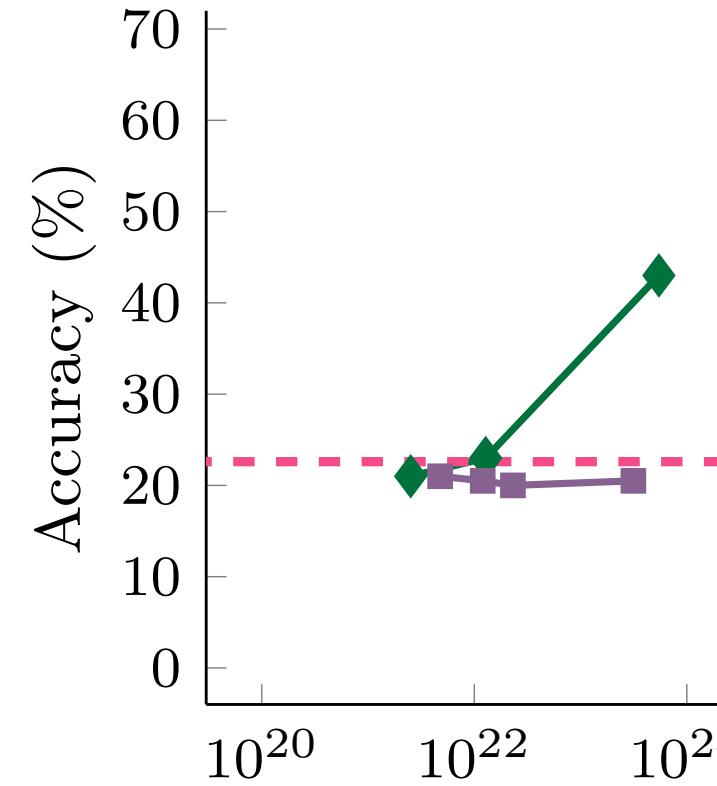
(C) Word unscramble



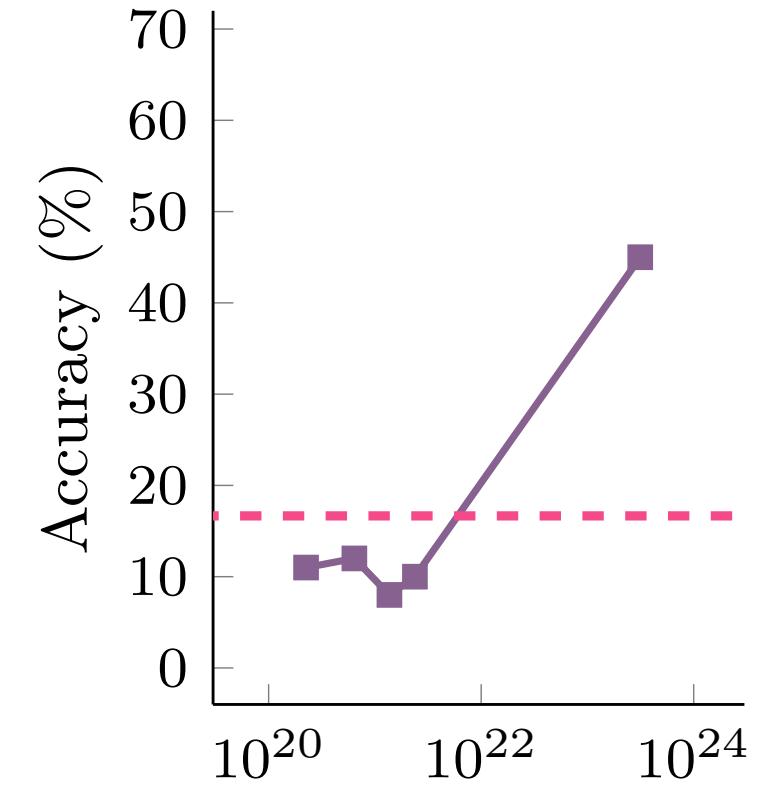
(D) Persian QA



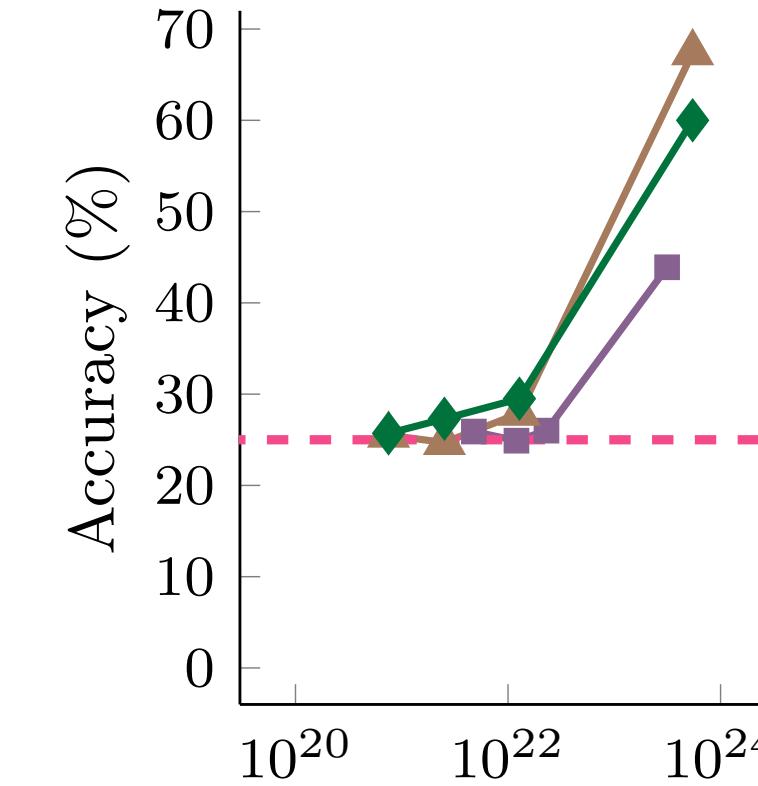
(E) TruthfulQA



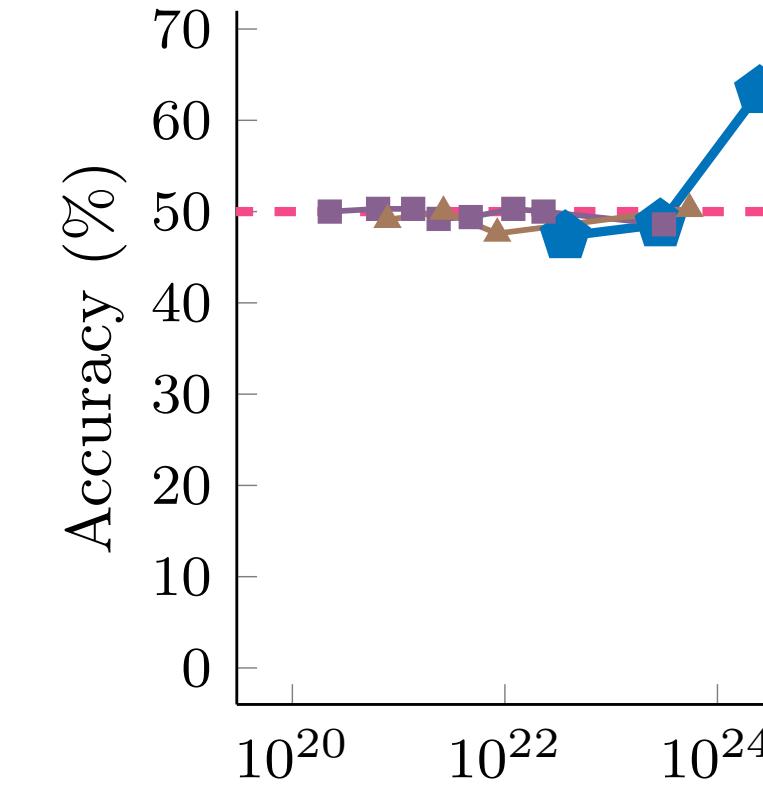
(F) Grounded mappings



(G) Multi-task NLU



(H) Word in context



Model scale (training FLOPs)

Are Emergent Abilities of Large Language Models a Mirage?

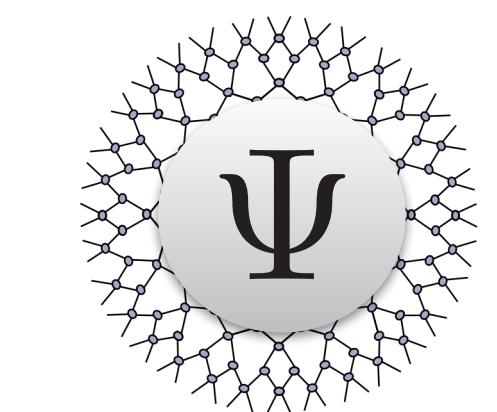
Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Computer Science, Stanford University

$$\text{Accuracy}(N) \approx p_N(\text{single token correct})^{\text{num. of tokens}} = \exp\left(-\left(N/c\right)^\alpha\right)^L$$

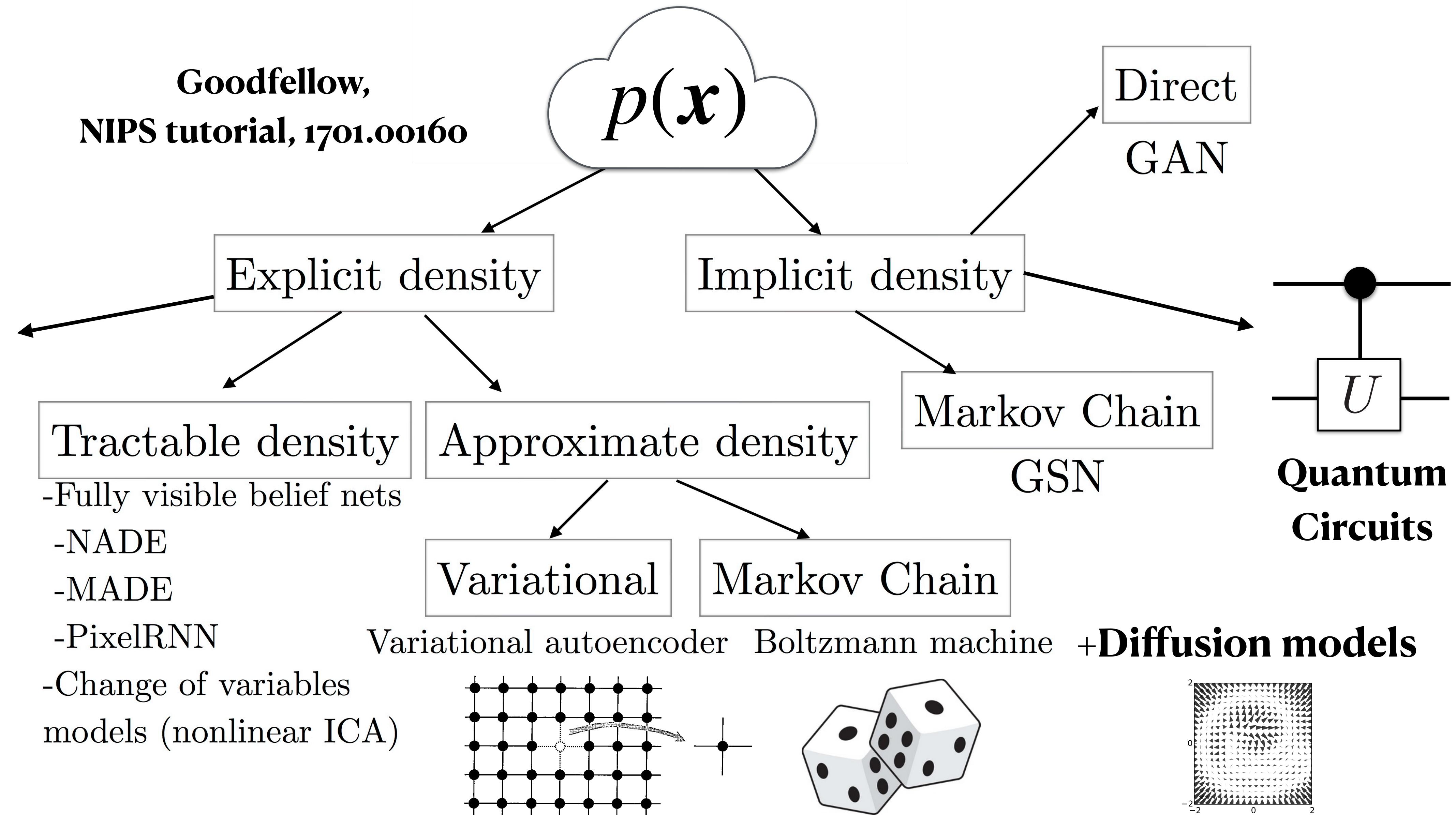
"The researcher's choice of metric can nonlinearly and/or discontinuously transform the error rate in a manner that causes the model performance to appear sharp and unpredictable."

Generative models and their physics genes



**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



Normalizing flows



Parallel WaveNet 1711.10433

<https://deepmind.com/blog/high-fidelity-speech-synthesis-wavenet/>



Glow 1807.03039

<https://blog.openai.com/glow/>

Normalizing flows



Parallel WaveNet 1711.10433

<https://deepmind.com/blog/high-fidelity-speech-synthesis-wavenet/>

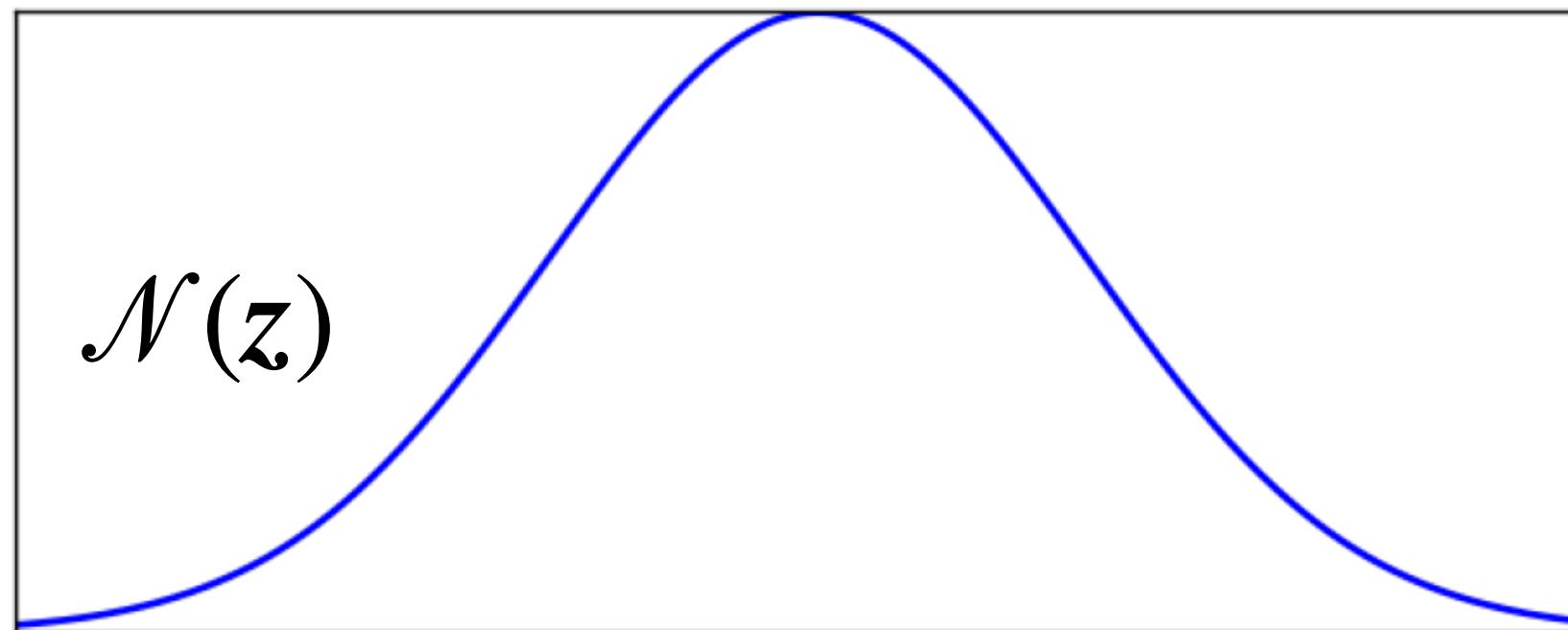


Glow 1807.03039

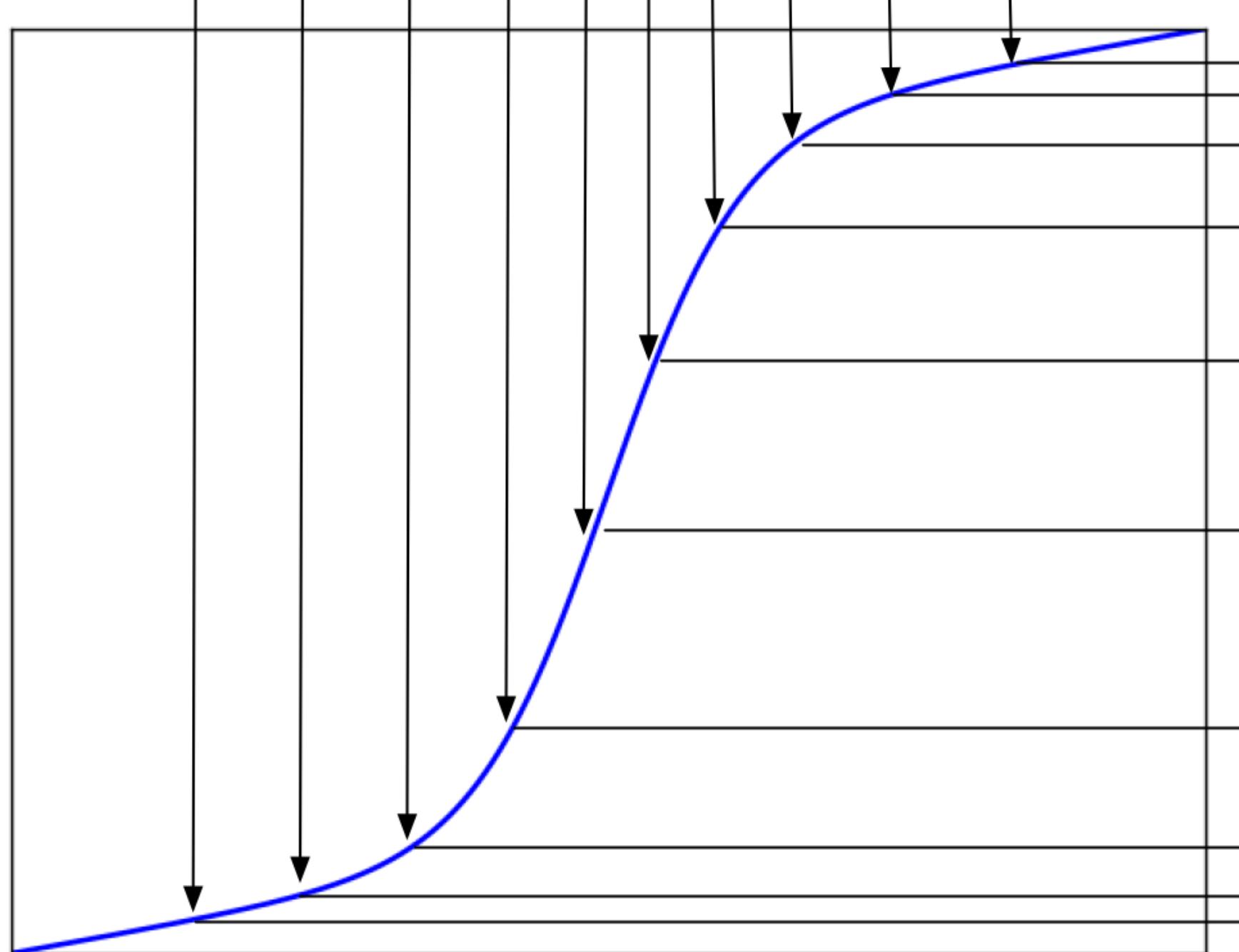
<https://blog.openai.com/glow/>

Normalizing flow in a nutshell

Base
density



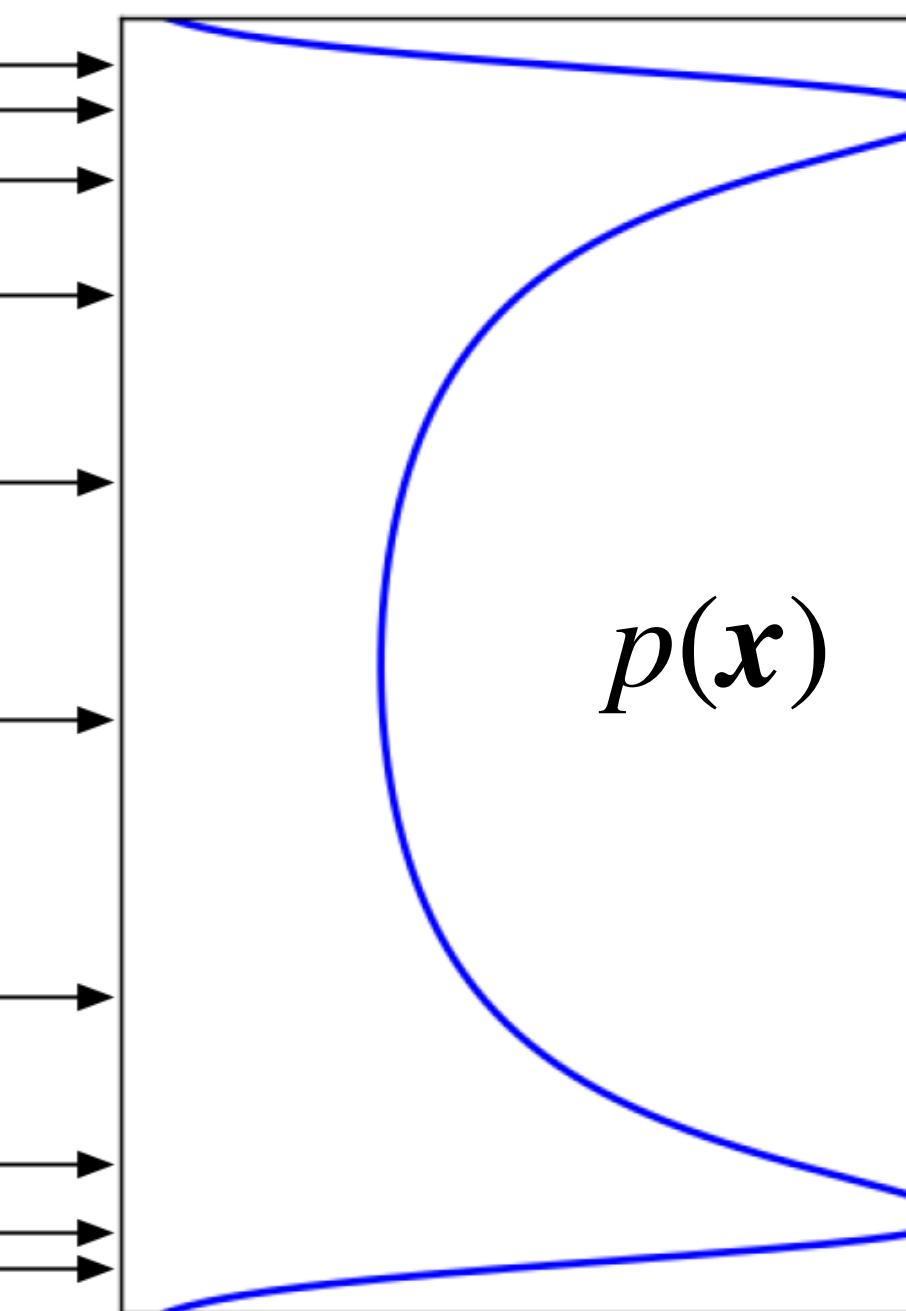
“neural net”
with 1 neuron



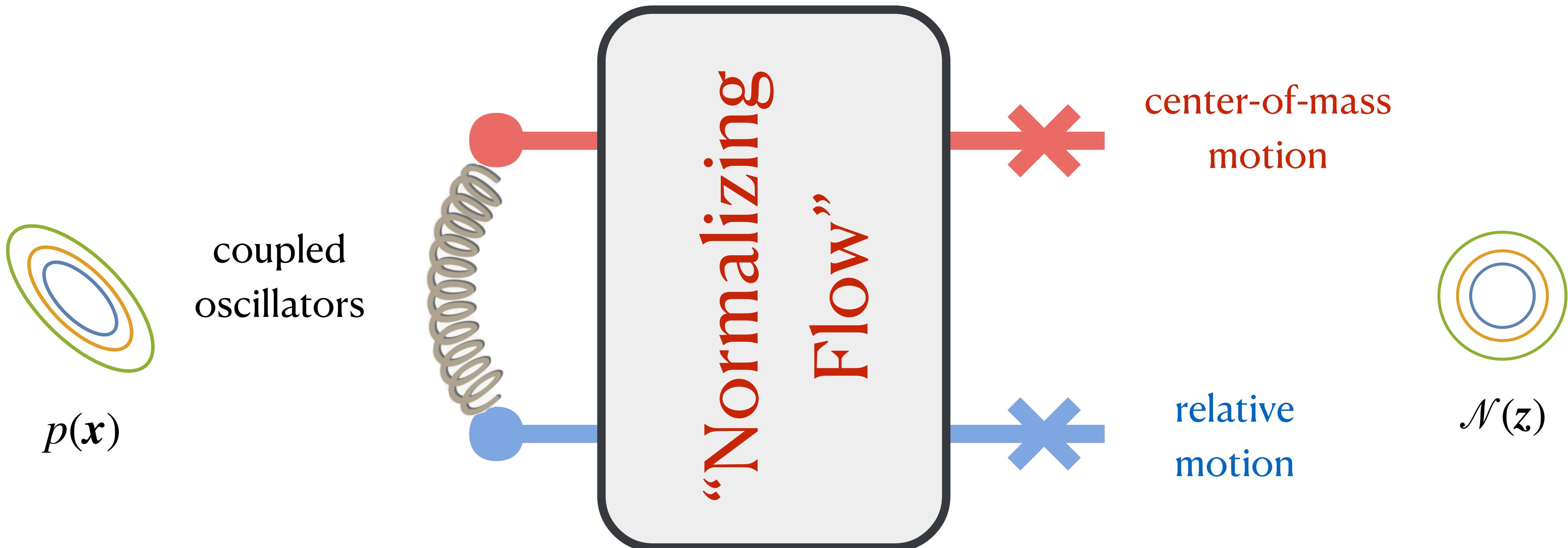
$$p(x) = \mathcal{N}(z) \left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$

Review article 1912.02762

Target
density



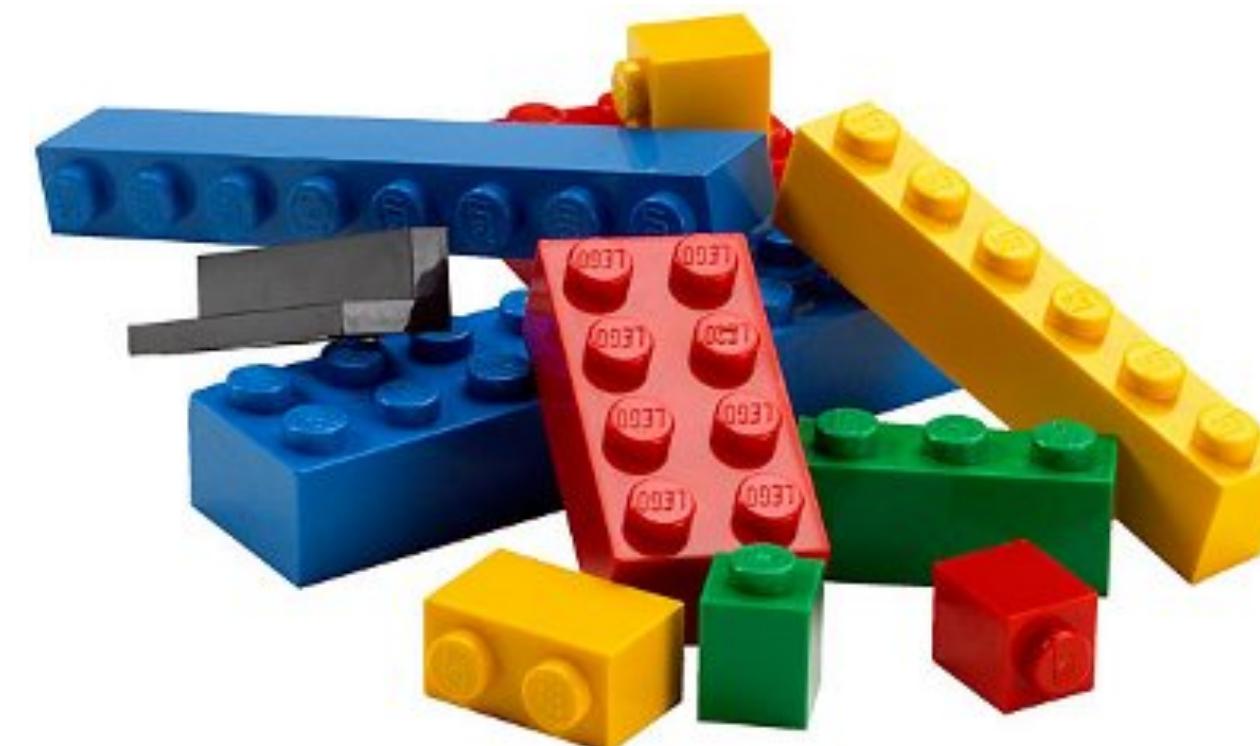
Physics intuition of normalizing flow



High-dimensional, nonlinear, learnable, composable transformations

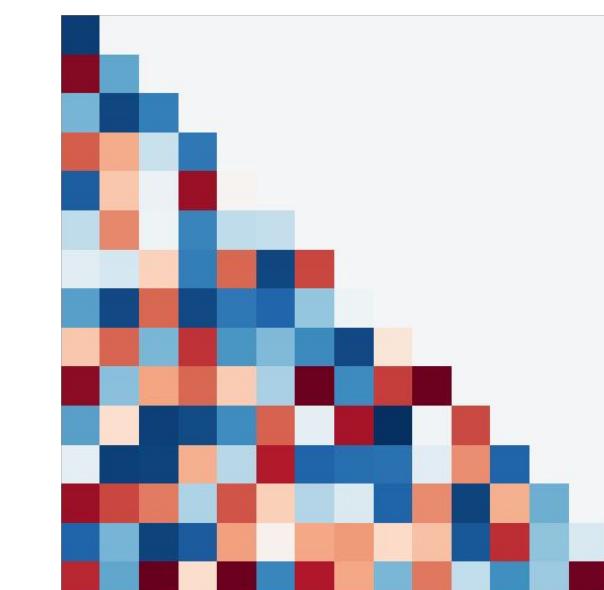
Flow architecture design

Composability

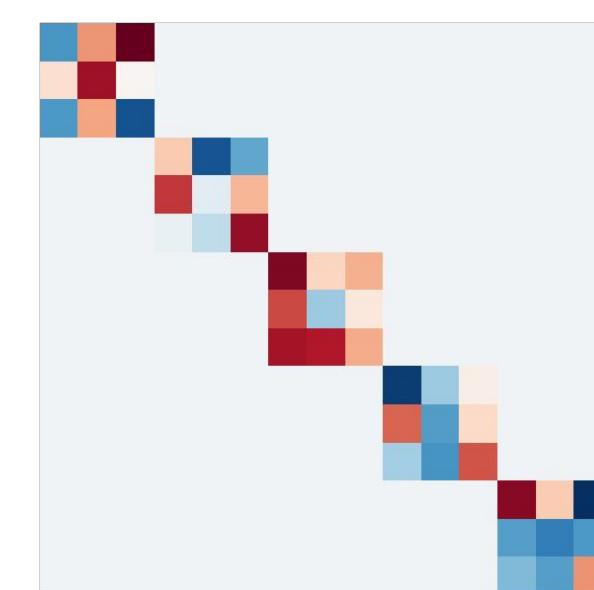


Balanced
efficiency &
inductive bias

$$\left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$



Autoregressive



Blockwise

$$z = \mathcal{T}(x)$$
$$\mathcal{T} = \mathcal{T}_1 \circ \mathcal{T}_2 \circ \mathcal{T}_3 \circ \dots$$

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t)v] = 0$$

Continuous flow

Example of a building block

Forward

$$\begin{cases} \mathbf{x}_< = \mathbf{z}_< \\ \mathbf{x}_> = \mathbf{z}_> \odot e^{s(\mathbf{z}_<)} + t(\mathbf{z}_<) \end{cases}$$

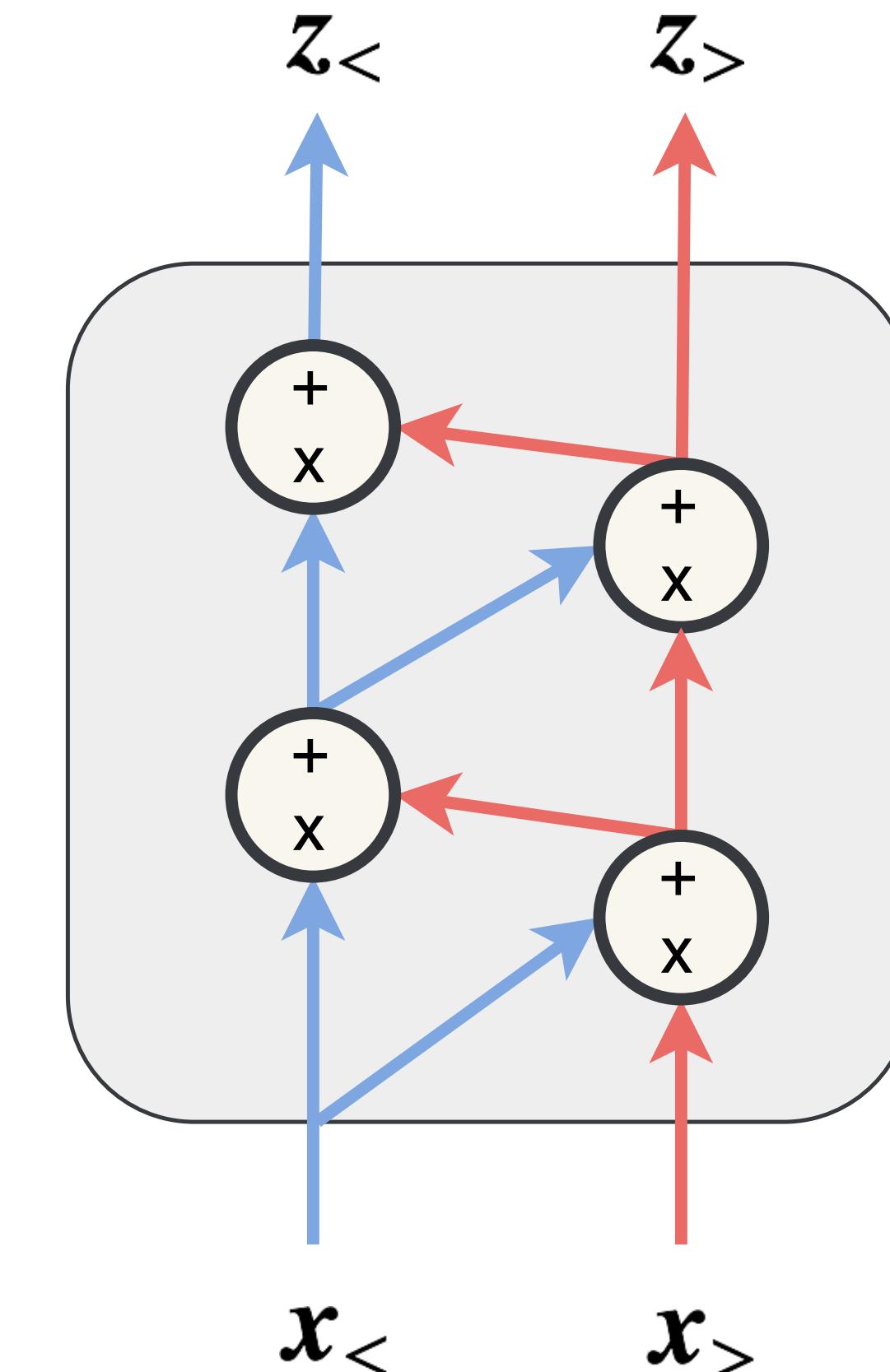
arbitrary
neural nets

Inverse

$$\begin{cases} \mathbf{z}_< = \mathbf{x}_< \\ \mathbf{z}_> = (\mathbf{x}_> - t(\mathbf{x}_<)) \odot e^{-s(\mathbf{x}_<)} \end{cases}$$

Log-Abs-Jacobian-Det

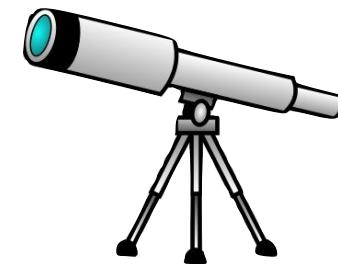
$$\ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i [s(\mathbf{z}_<)]_i$$



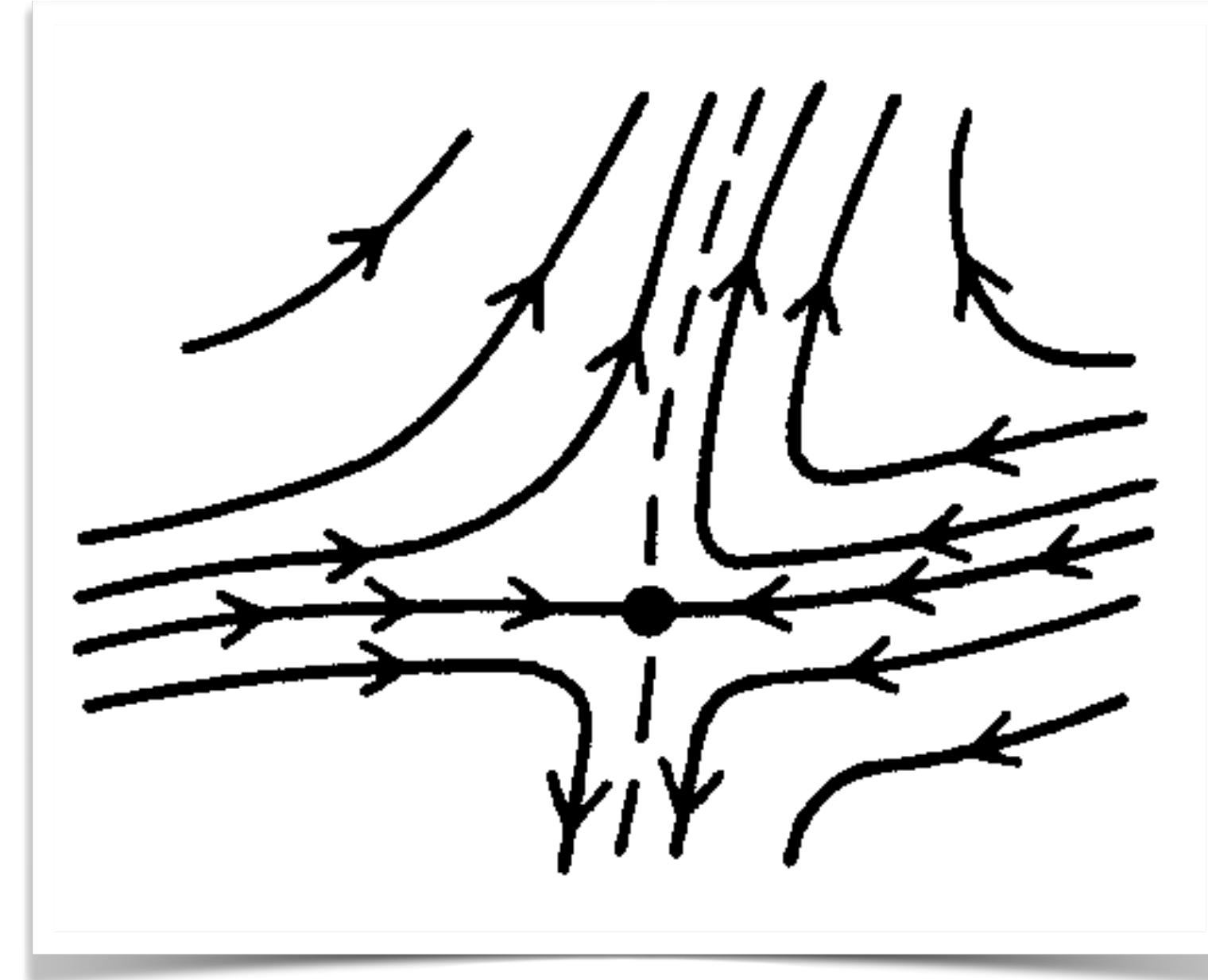
Real NVP, Dinh et al, 1605.08803

Turns out to have surprising connection Störmer–Verlet integration

Why is flow useful for physics?



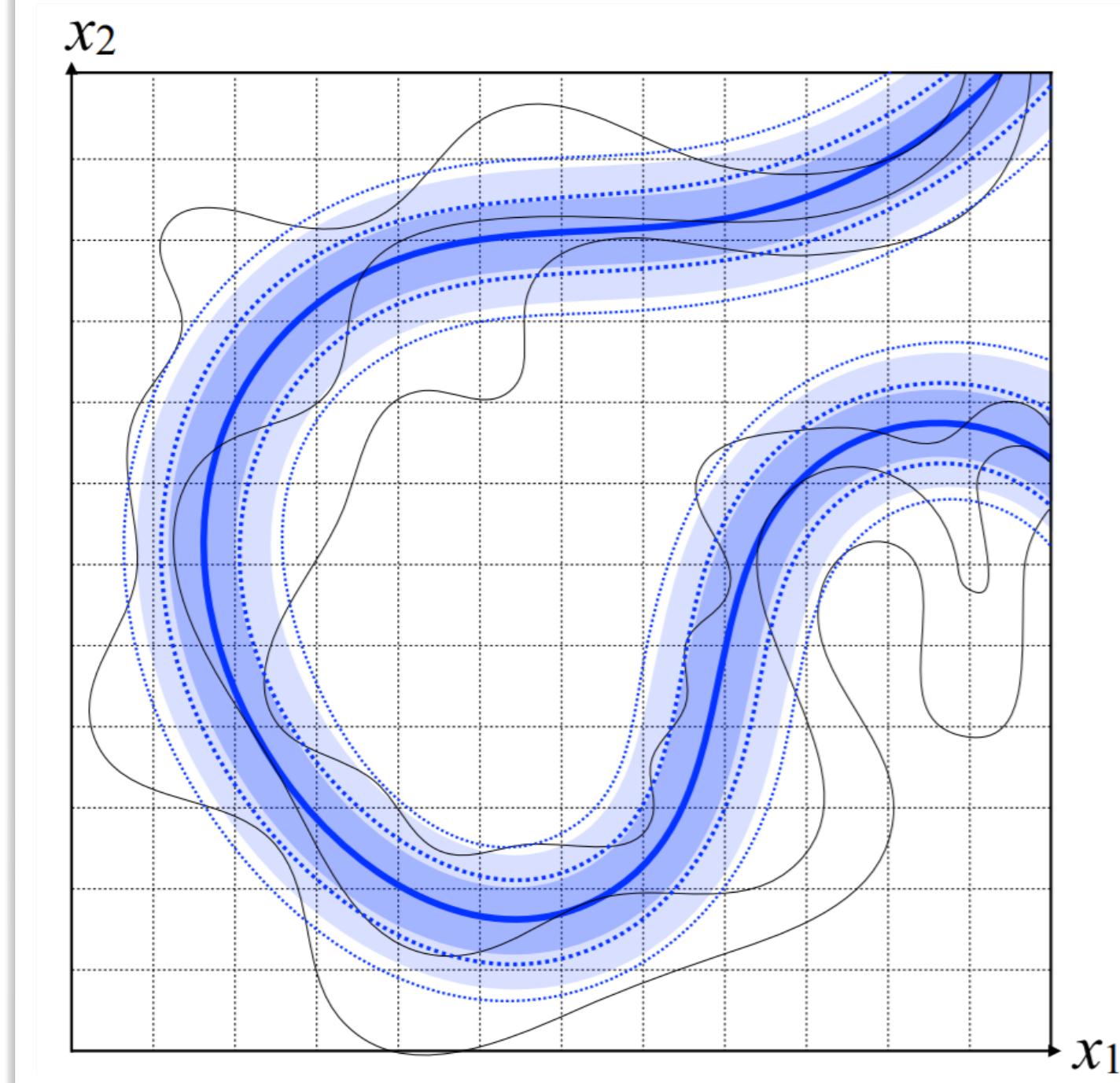
Renormalization group



Effective theory emerges upon
transformation of the variables



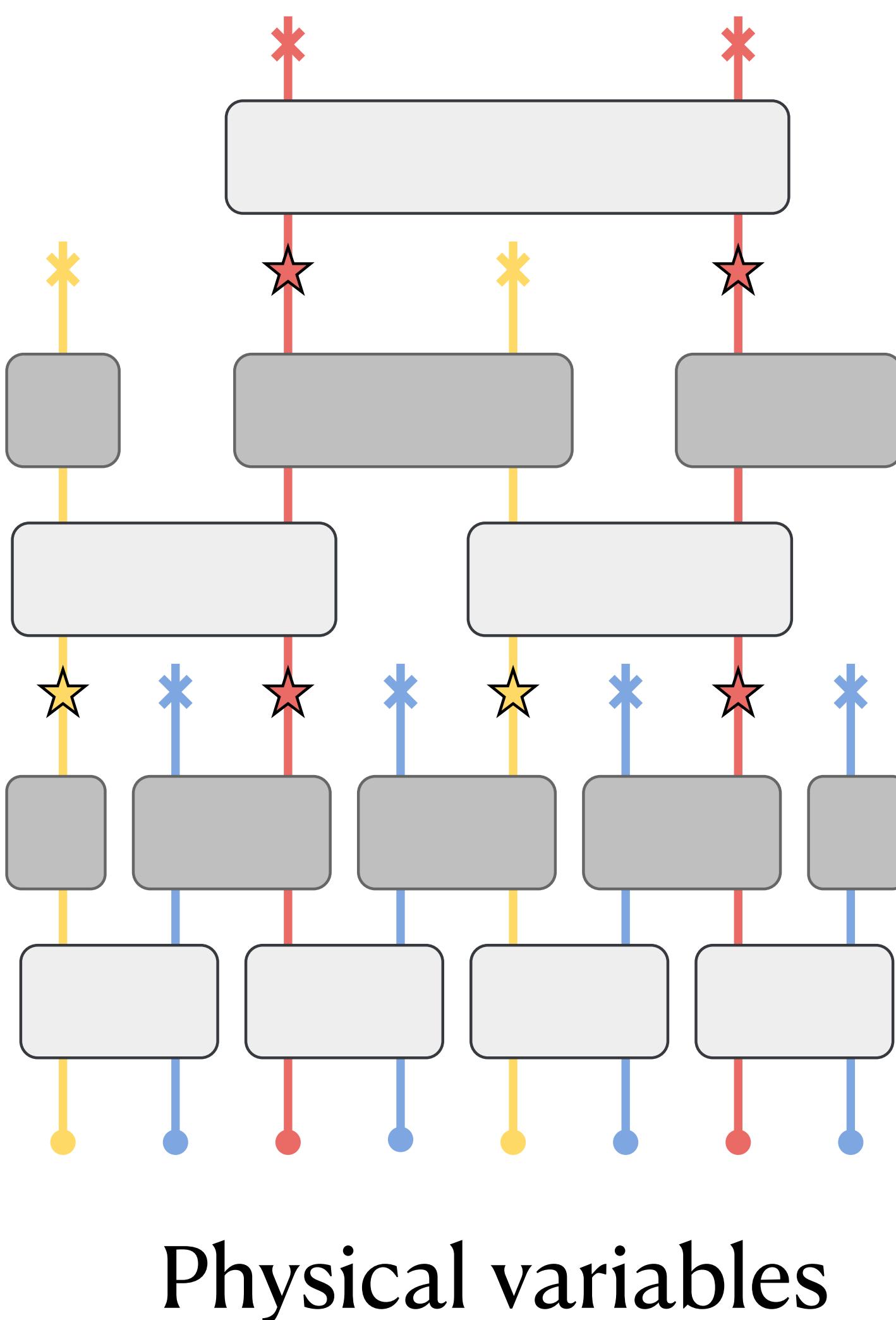
Monte Carlo update



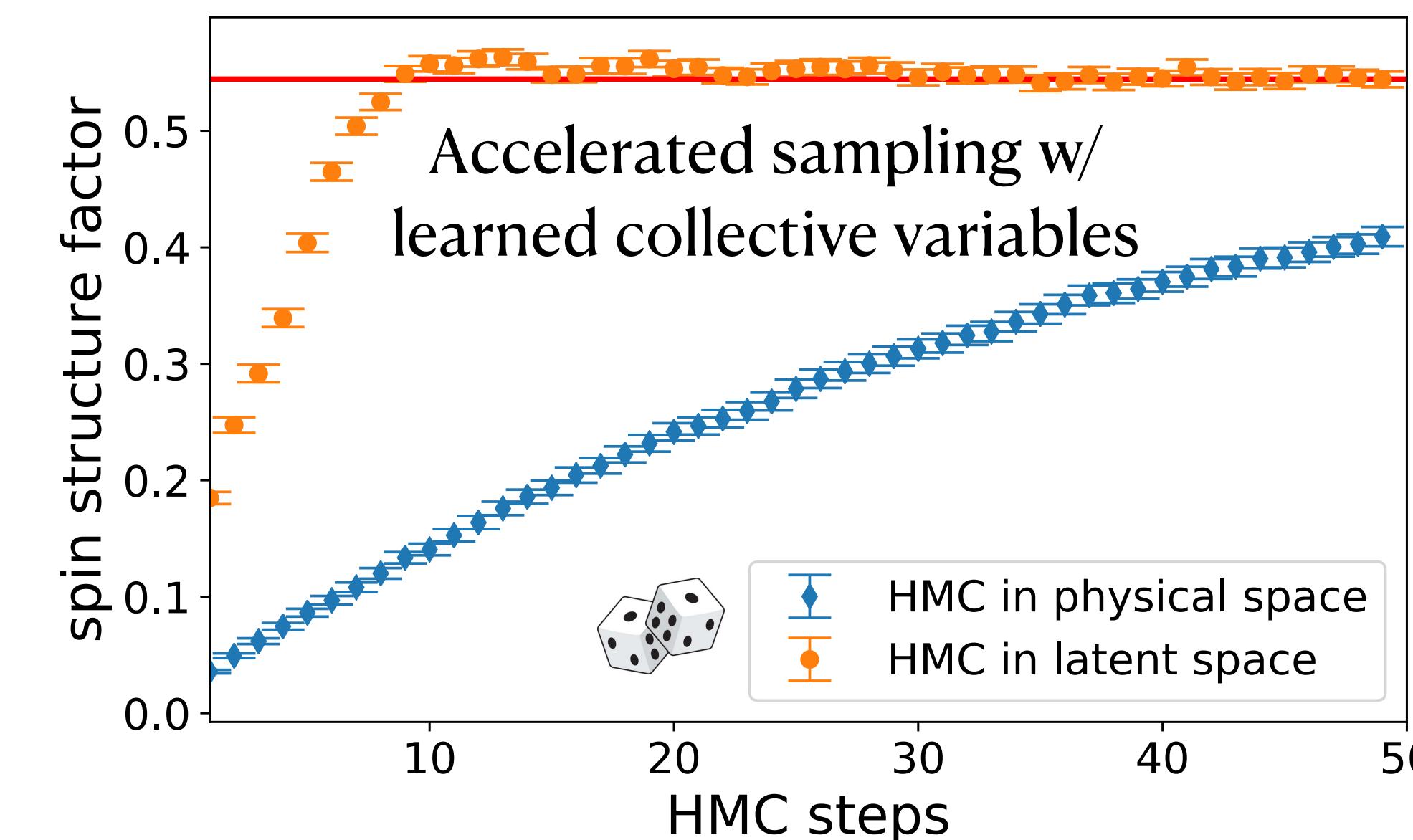
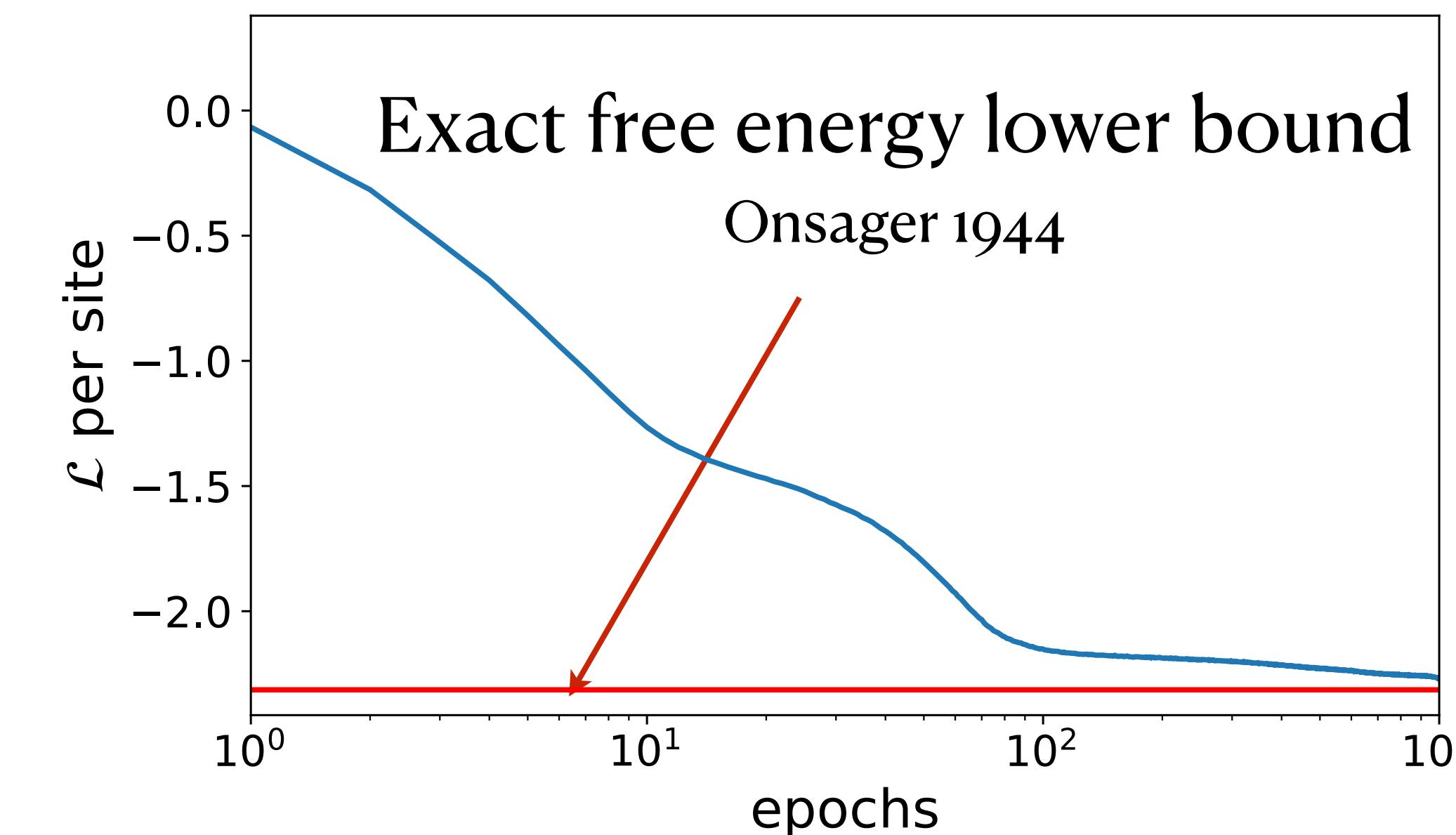
Physics happens on a manifold
Train neural nets to unfold that manifold

Neural network renormalization group

Collective variables

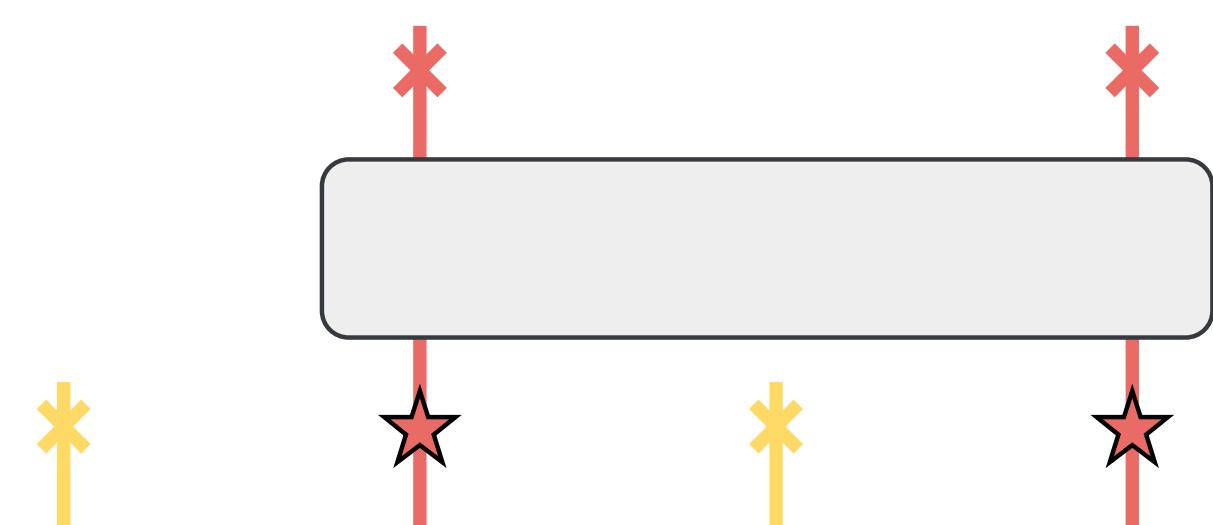


Li, LW, PRL '18 [lio12589/NeuralRG](#)



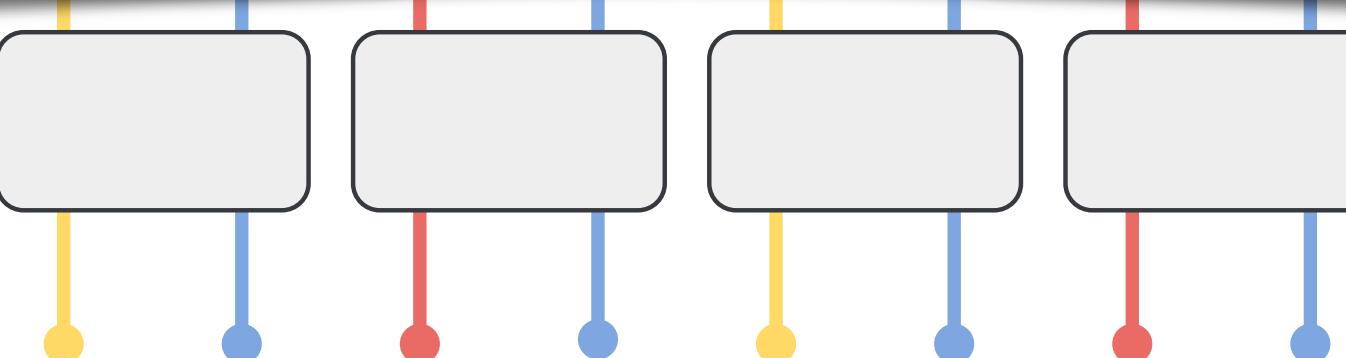
Neural network renormalization group

Collective variables



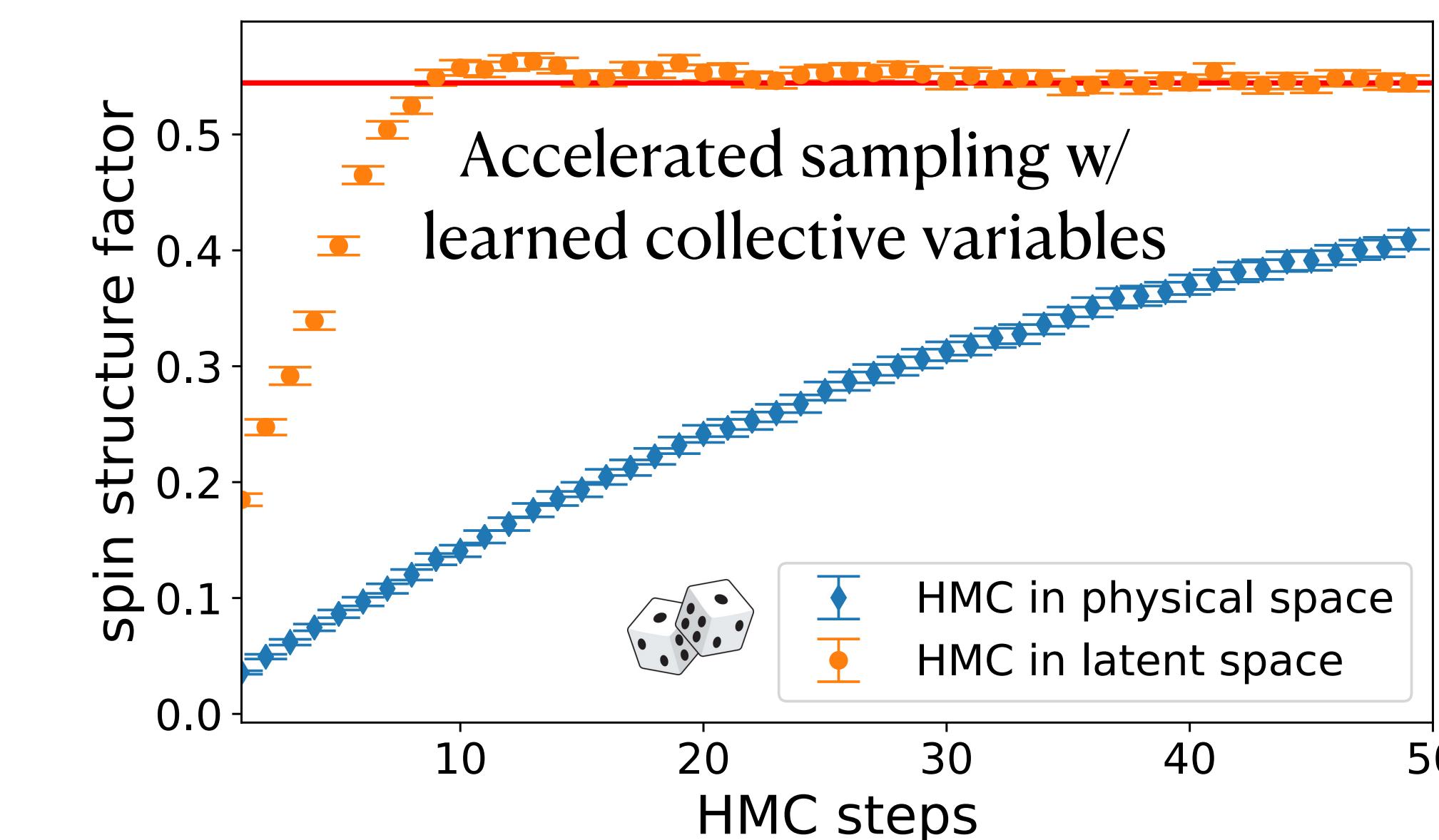
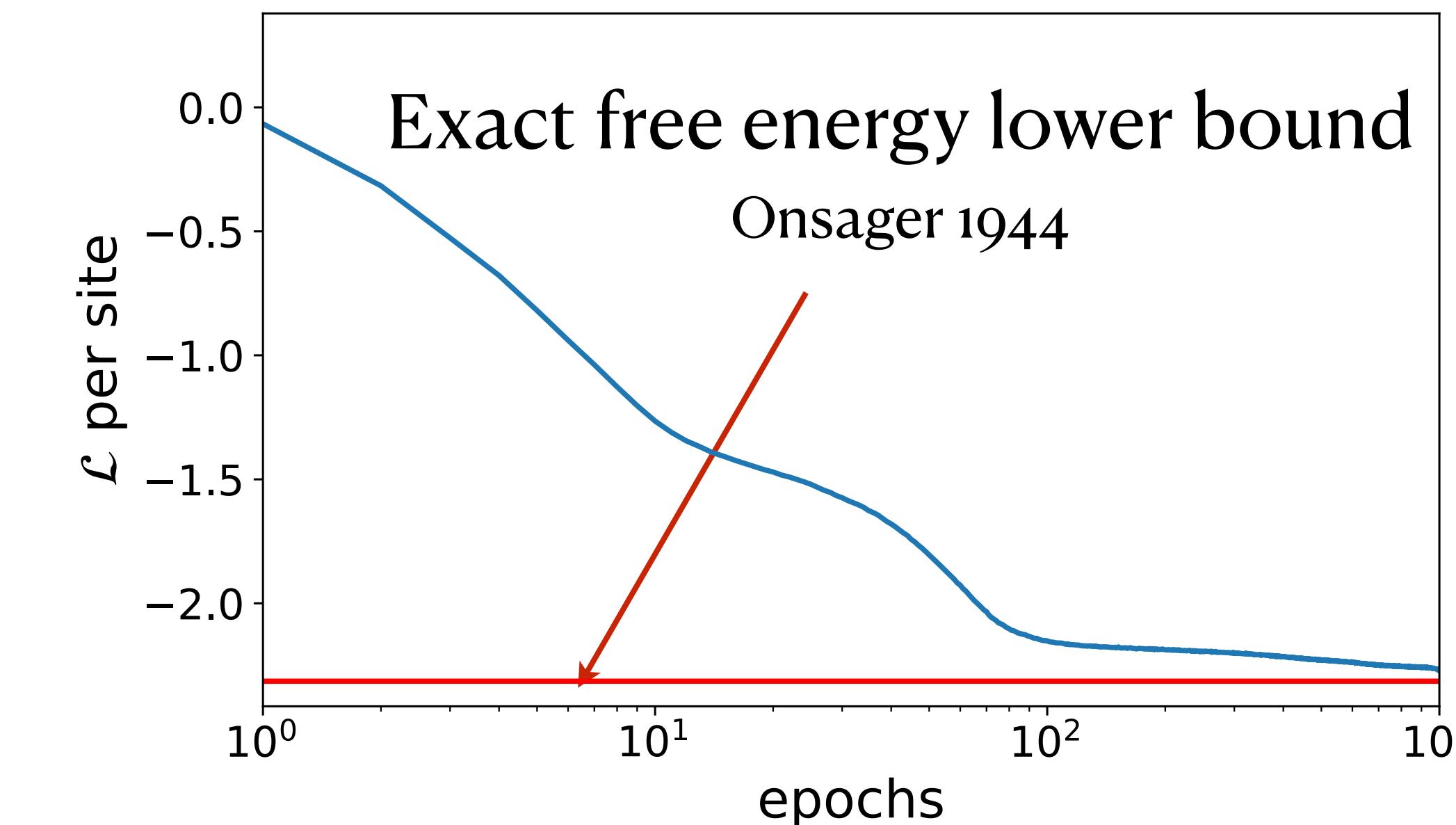
Probability Transformation

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(\mathbf{z}) - \ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right|$$

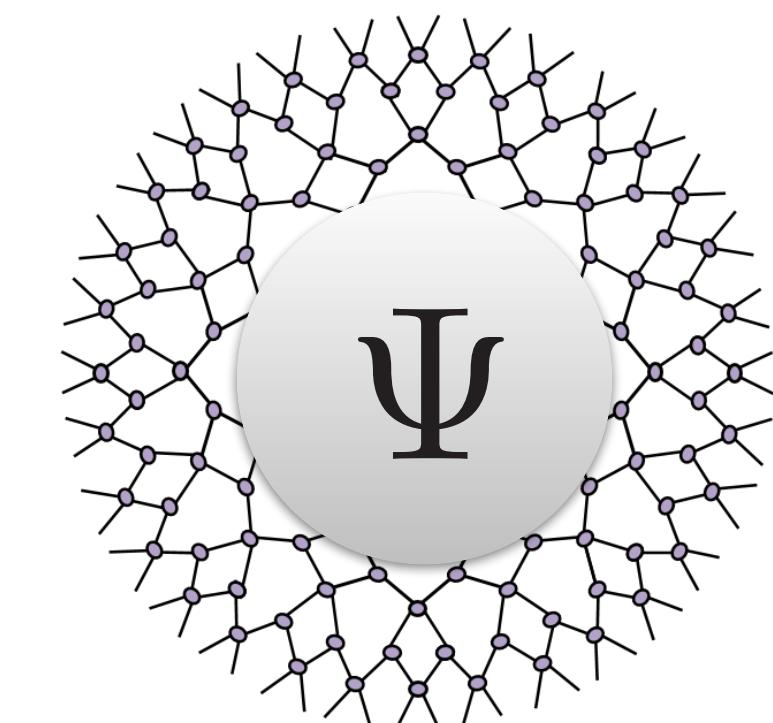
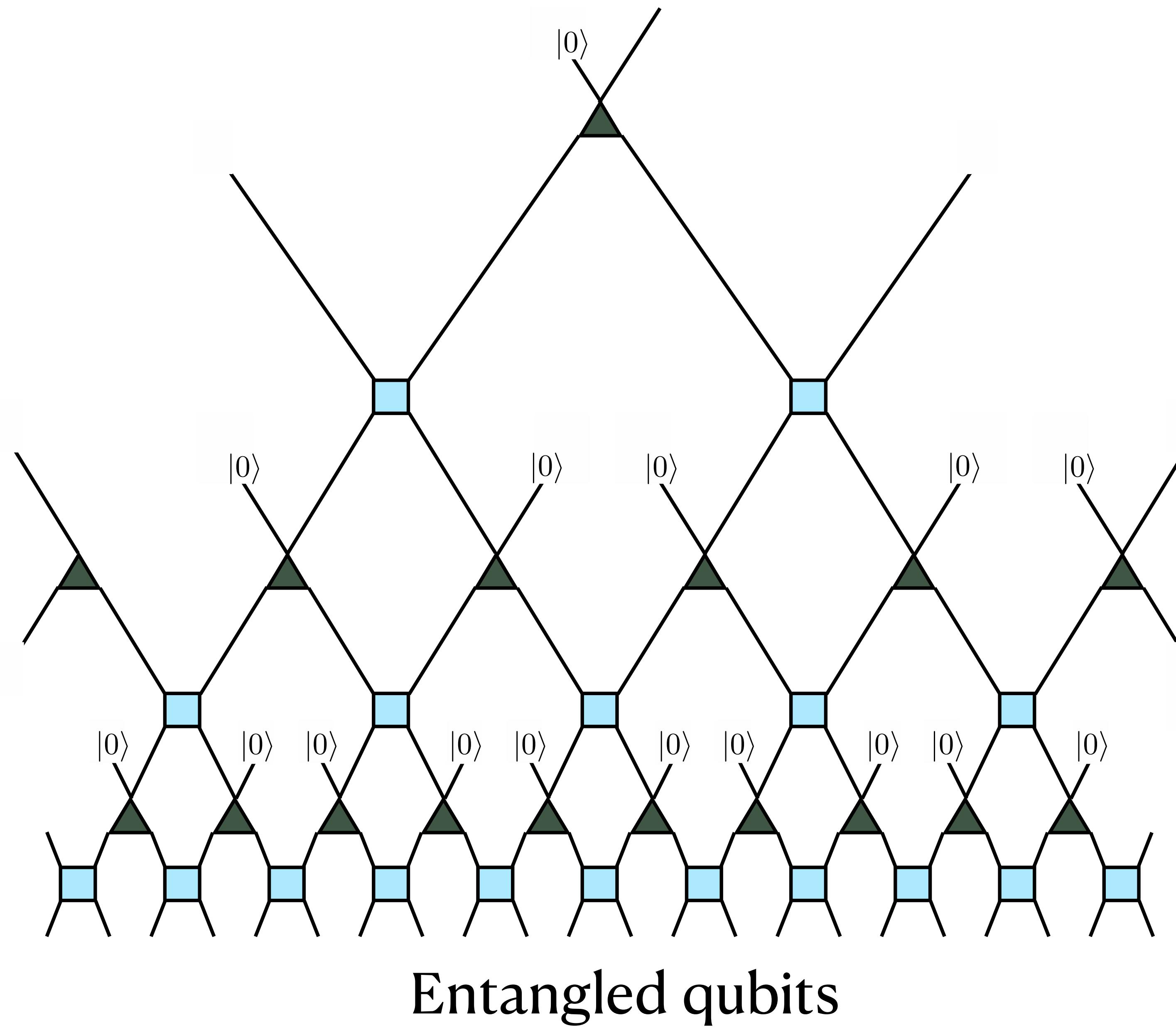


Physical variables

Li, LW, PRL '18 [lio12589/NeuralRG](#)

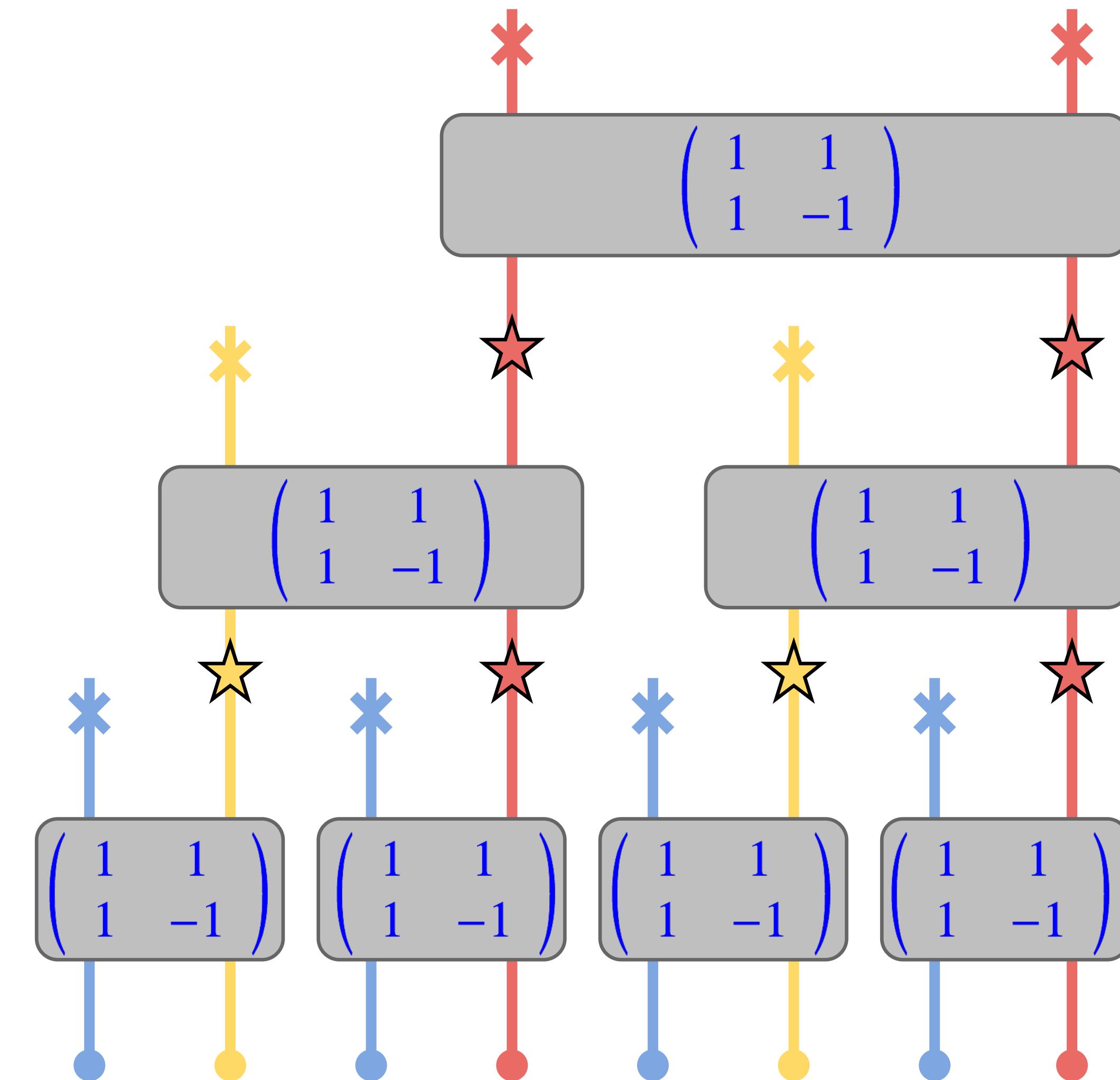
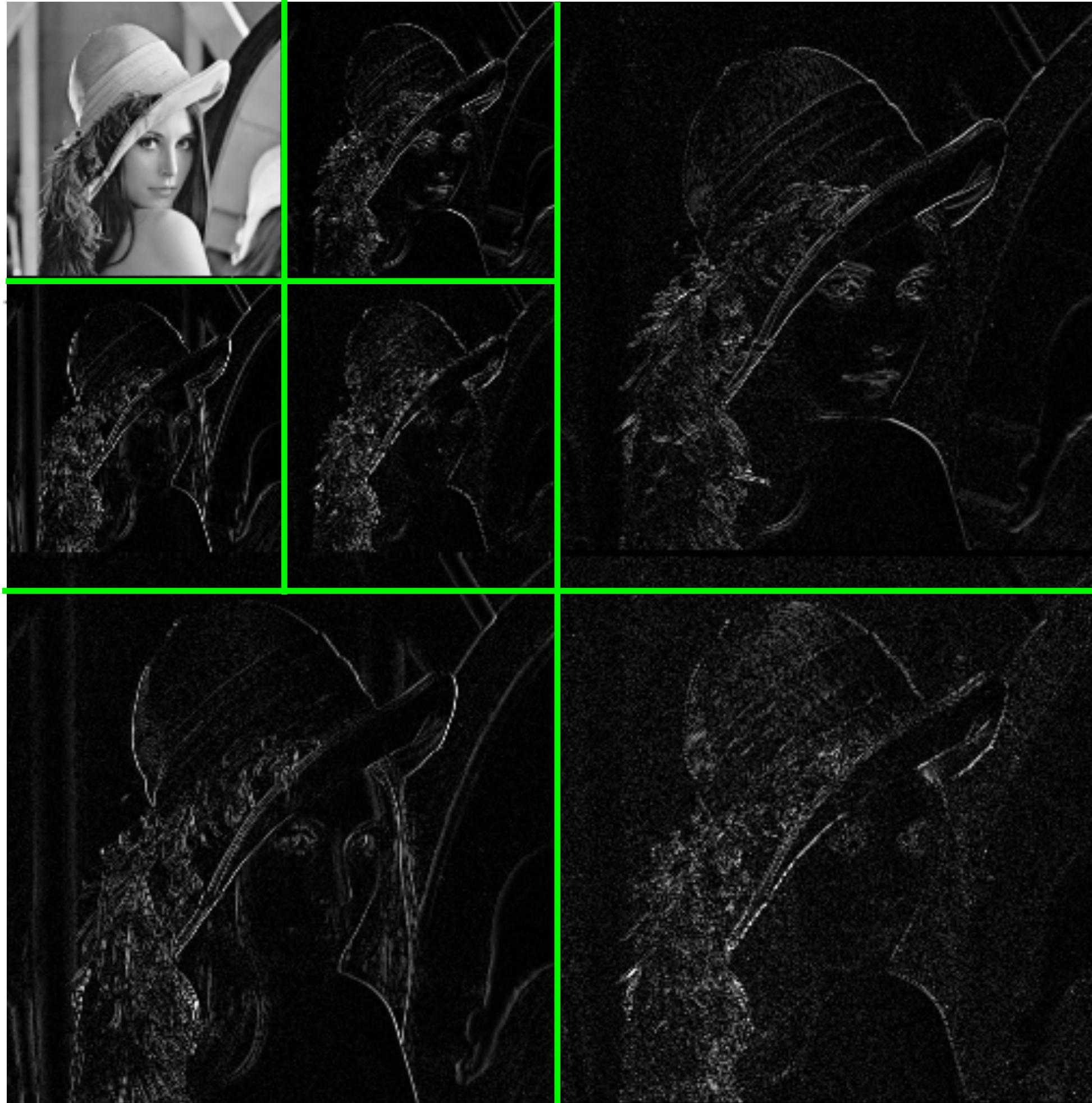


Quantum version of the architecture



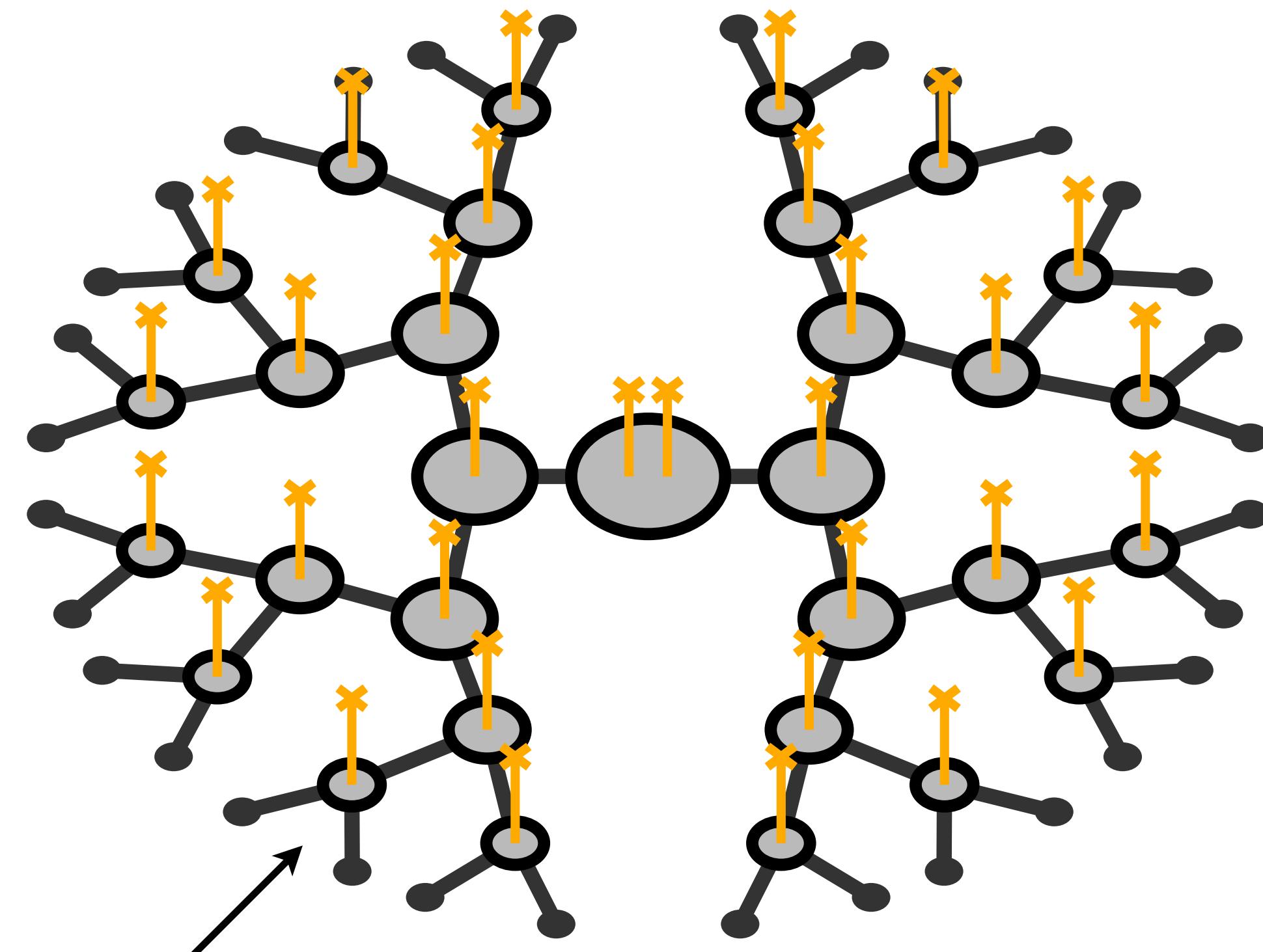
**Multi-Scale
Entanglement
Renormalization
Ansatz**

Connection to wavelets



Nonlinear & adaptive generalizations of wavelets

Neural network holographic RG



Bijective
neural net

Physical variables on the boundary

Latent variables in the bulk

RG flows along the radial direction

Information is preserved by the flow

Hu, Li, LW, You, PRR '20

See also Hashimoto et al 1809.10536, 2006.00712

Mutual information reveals the emergent geometry in the bulk

Continuous normalizing flows

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(\mathbf{z}) - \ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right|$$

Consider infinitesimal change-of-variables Chen et al 1806.07366

$$\mathbf{x} = \mathbf{z} + \varepsilon \boldsymbol{\nu}$$

$$\ln p(\mathbf{x}) - \ln \mathcal{N}(\mathbf{z}) = - \ln \left| \det \left(1 + \varepsilon \frac{\partial \boldsymbol{\nu}}{\partial \mathbf{z}} \right) \right|$$

$$\varepsilon \rightarrow 0$$

$$\frac{d\mathbf{x}}{dt} = \boldsymbol{\nu}$$

$$\frac{d \ln p(\mathbf{x}, t)}{dt} = - \nabla \cdot \boldsymbol{\nu}$$

Continuous normalizing flows

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(\mathbf{z}) - \ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right|$$

Consider infinitesimal change-of-variables Chen et al 1806.07366

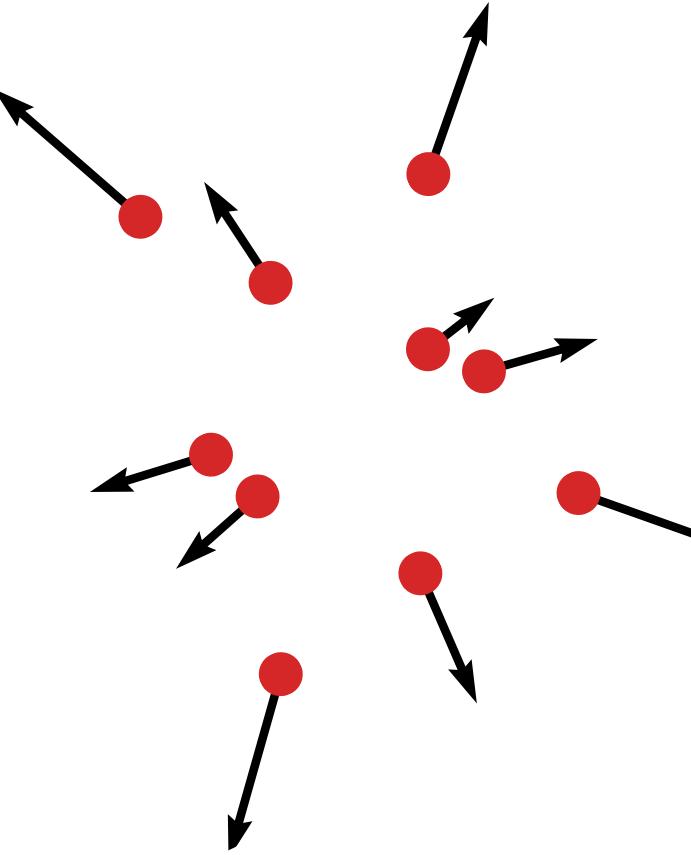
$$\begin{aligned}\mathbf{x} &= \mathbf{z} + \varepsilon \boldsymbol{\nu} \\ \varepsilon &\rightarrow 0\end{aligned}$$

$$\frac{d\mathbf{x}}{dt} = \boldsymbol{\nu}$$

$$\begin{aligned}\ln p(\mathbf{x}) - \ln \mathcal{N}(\mathbf{z}) &= - \ln \left| \det \left(1 + \varepsilon \frac{\partial \boldsymbol{\nu}}{\partial \mathbf{z}} \right) \right| \\ t = 0 &\quad t = 1 \\ \frac{d \ln p(\mathbf{x}, t)}{dt} &= - \nabla \cdot \boldsymbol{\nu}\end{aligned}$$

Fluid physics behind flows

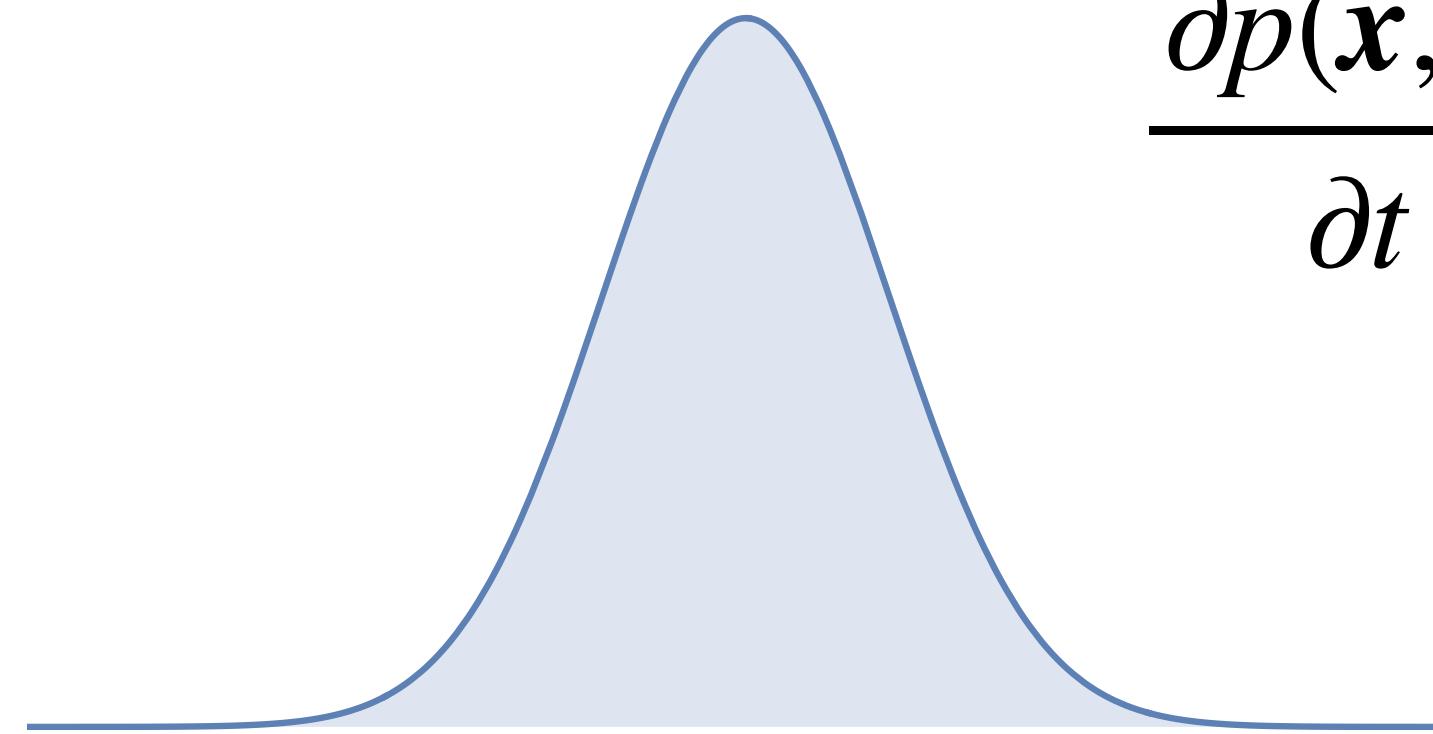
$$\frac{dx}{dt} = v$$
$$\frac{d \ln p(x, t)}{dt} = - \nabla \cdot v$$



Zhang, E, LW 1809.10188
[wangleiphy/MongeAmpereFlow](#)

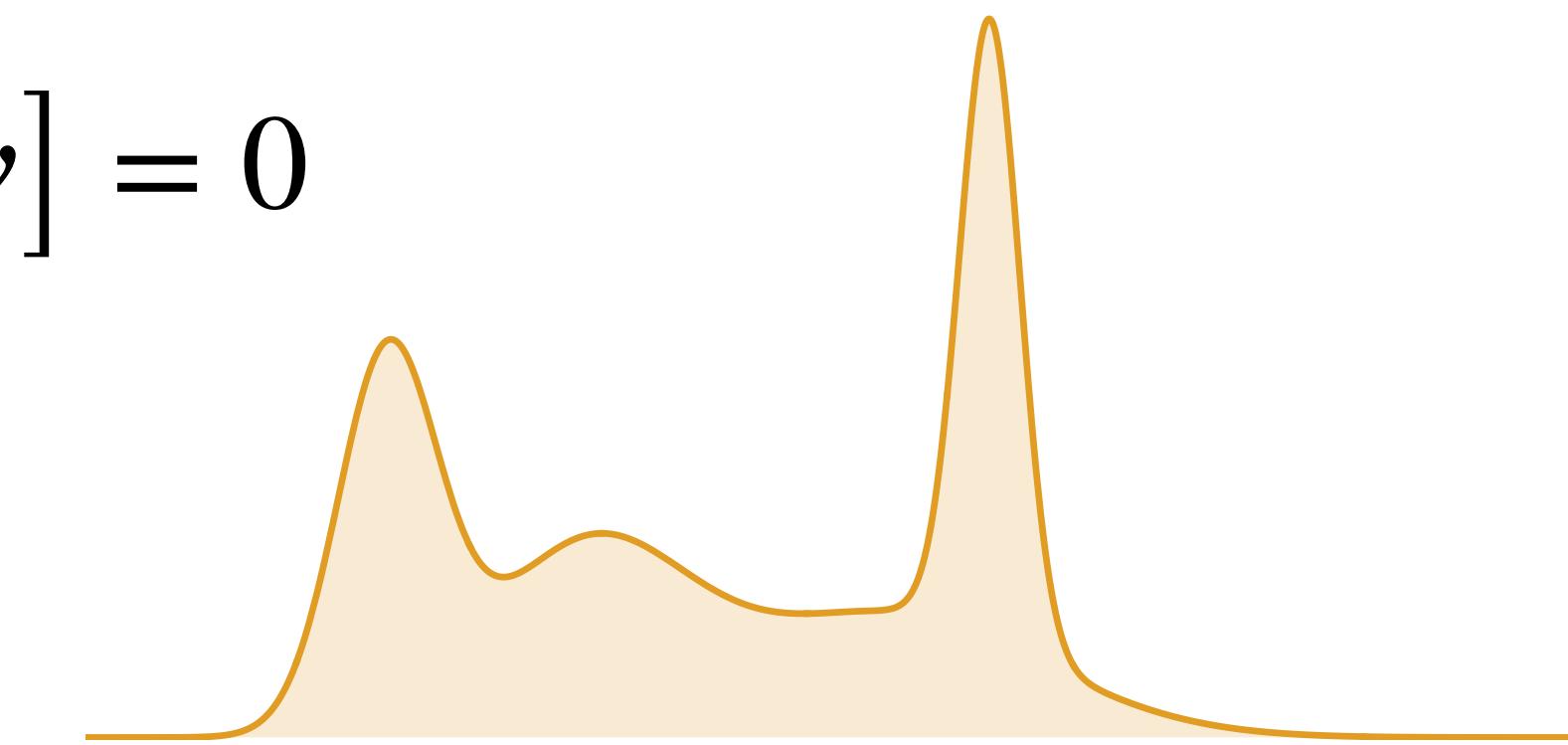
$$\frac{d}{dt} = \frac{\partial}{\partial t} + v \cdot \nabla \quad \text{"material derivative"}$$

Lagrangian v.s. Euler approach to fluid mechanics



Simple density

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t)v] = 0$$



Complex density

Infinitesimal Flows Another way to reduce the computational overhead of normalizing flows is to use an ordinary differential equation to generate f [2]. In this case, the probability distribution changes over a finite time from $p(\mathbf{x}; 0)$ to $p(\mathbf{z}; T)$, where \mathbf{z} is the end point of a curve defined by the ODE $\dot{\mathbf{x}}(t) = \mathbf{v}(\mathbf{x}(t))$, $\mathbf{x}(0) = \mathbf{x}$. For a small time step dt , we can approximate $\mathbf{x}(t + dt)$ to first order as $\mathbf{x}(t + dt) = \mathbf{x}(t) + dt\mathbf{v}(\mathbf{x}(t)) + \mathcal{O}(dt^2)$. Plugging this into Eq. 2 yields:

$$\log p(\mathbf{x} + dt\mathbf{v}(\mathbf{x}) + \mathcal{O}(dt^2); t + dt) = \log p(\mathbf{x}; t) - \log |\mathbf{J}_f(\mathbf{x})| \quad (3)$$

$$= \log p(\mathbf{x}; t) - \log |\mathbf{I} + dt\mathbf{J}_v(\mathbf{x}) + \mathcal{O}(dt^2)| \quad (4)$$

Taking a Taylor series gives:

$$\log p(\mathbf{x}; t + dt) + dt\mathbf{v}(\mathbf{x})^T \nabla \log p(\mathbf{x}; t + dt) = \log p(\mathbf{x}; t) - dt \text{Tr}(\mathbf{J}_v(\mathbf{x})) + \mathcal{O}(dt^2) \quad (5)$$

which, in the limit as $dt \rightarrow 0$, becomes:

$$\frac{\partial \log p(\mathbf{x}; t)}{\partial t} = -\mathbf{v}(\mathbf{x})^T \nabla \log p(\mathbf{x}; t) - \text{Tr}(\mathbf{J}_v(\mathbf{x})) = -\mathbf{v}^T \nabla \log p(\mathbf{x}; t) - \nabla \cdot \mathbf{v} \quad (6)$$

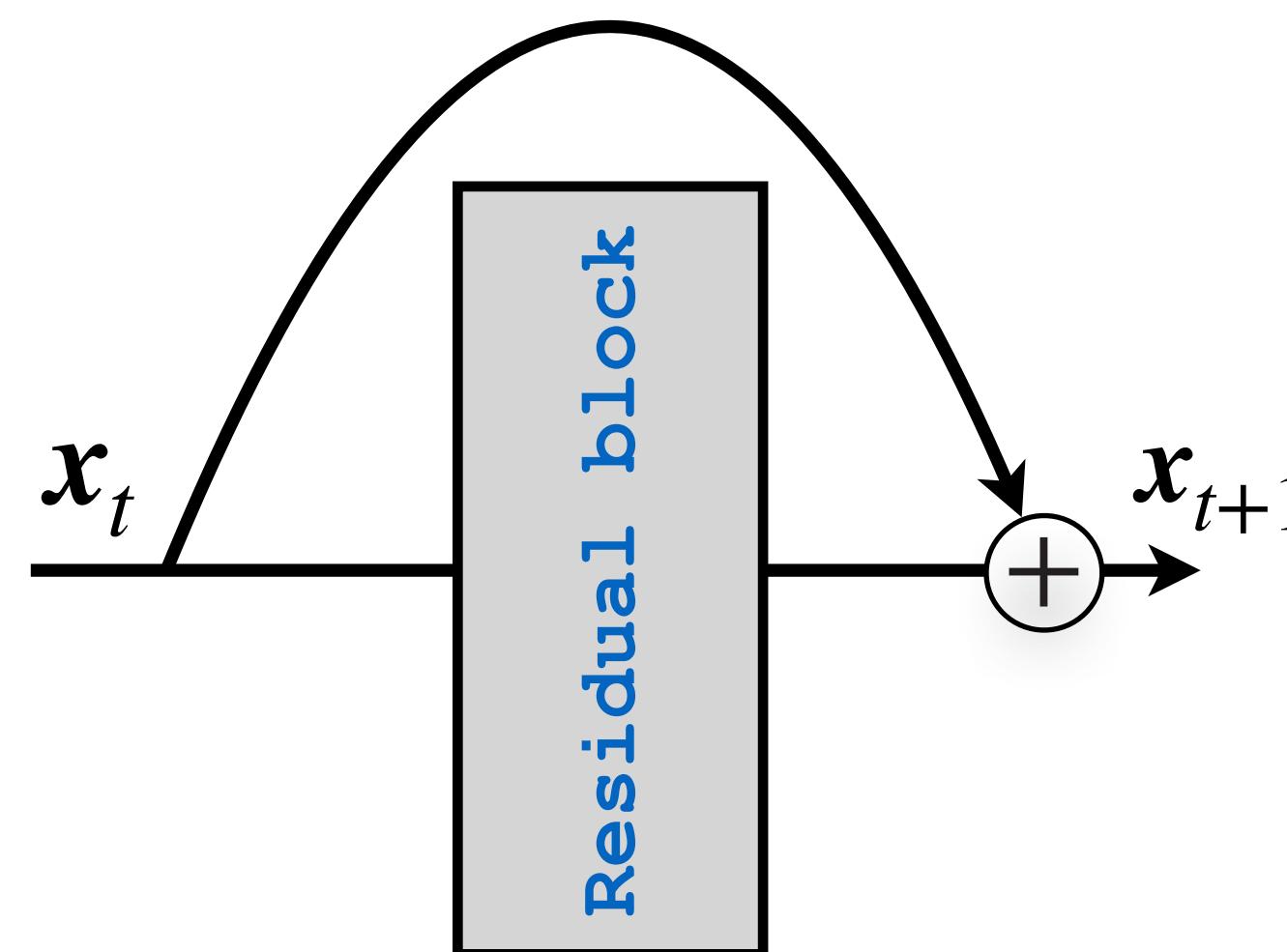
after some rearranging of terms. Here $\nabla \cdot$ is the divergence of a vector field, which is just another way of writing the trace of the Jacobian. The right-hand side of this equation is also the trace of the *Stein operator* of the distribution $p(\mathbf{x})$ applied to the function $\mathbf{v}(\mathbf{x})$, and plays a critical role in Stein variational gradient descent (SVGD) [13]. Switching from the log density to the density (and dropping the t for clarity), we find this expression can be simplified considerably:

$$\begin{aligned} \frac{1}{p(\mathbf{x})} \frac{\partial p(\mathbf{x})}{\partial t} &= -\mathbf{v}^T \frac{\nabla p(\mathbf{x})}{p(\mathbf{x})} - \nabla \cdot \mathbf{v} \\ \frac{\partial p(\mathbf{x})}{\partial t} &= -\mathbf{v}^T \nabla p(\mathbf{x}) - p(\mathbf{x}) \nabla \cdot \mathbf{v} \\ &= -\nabla \cdot (\mathbf{v}(\mathbf{x}) p(\mathbf{x})) \end{aligned} \quad (7)$$

This may also be familiar as the drift term of the Fokker-Planck equation [11, Eq. 6.48] or the continuity equation for conservation of mass in fluid mechanics. We will denote the change to a

Neural Ordinary Differential Equations

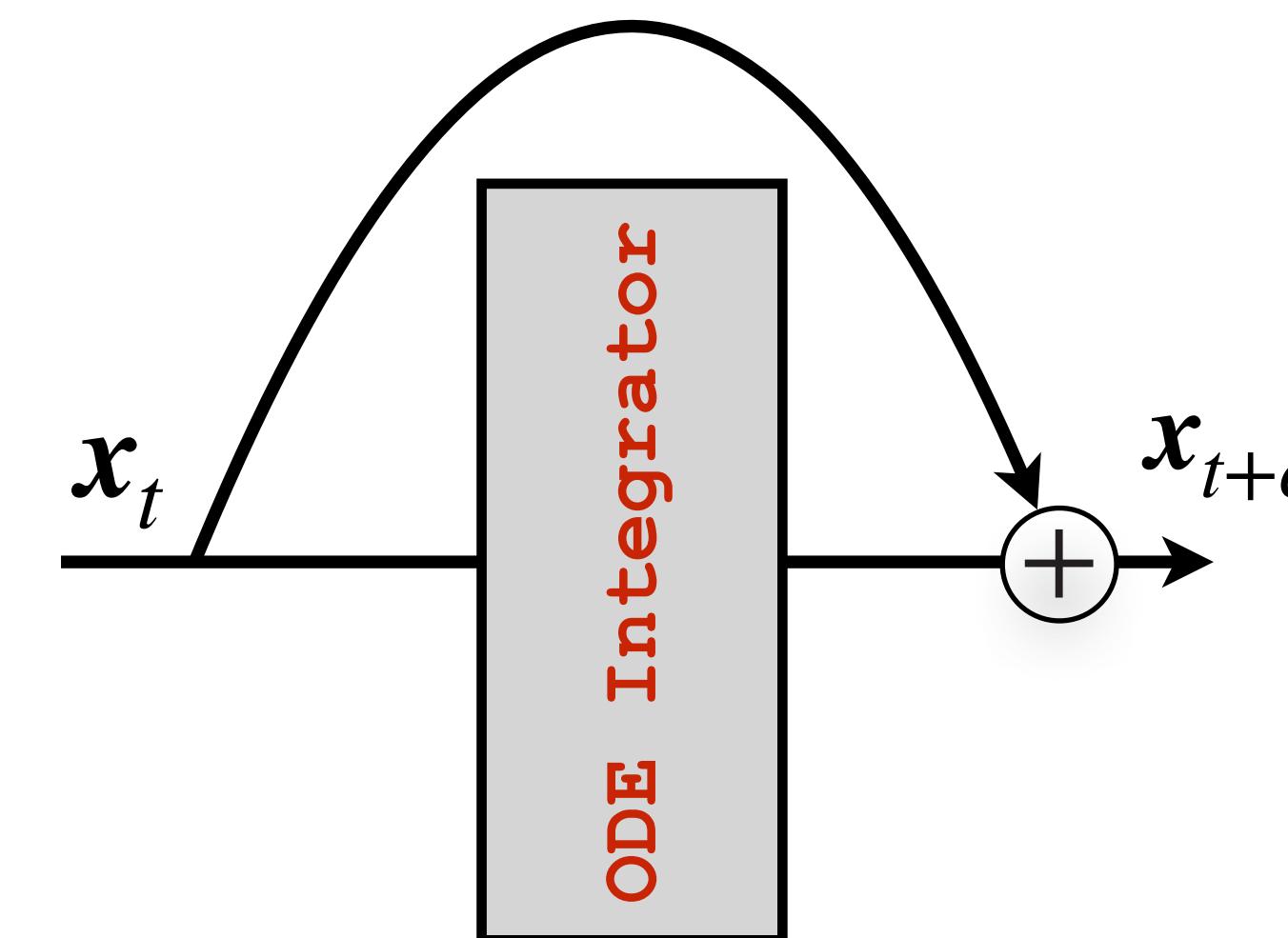
Residual network



$$x_{t+1} = x_t + v(x_t)$$

Chen et al, 1806.07366

ODE integration

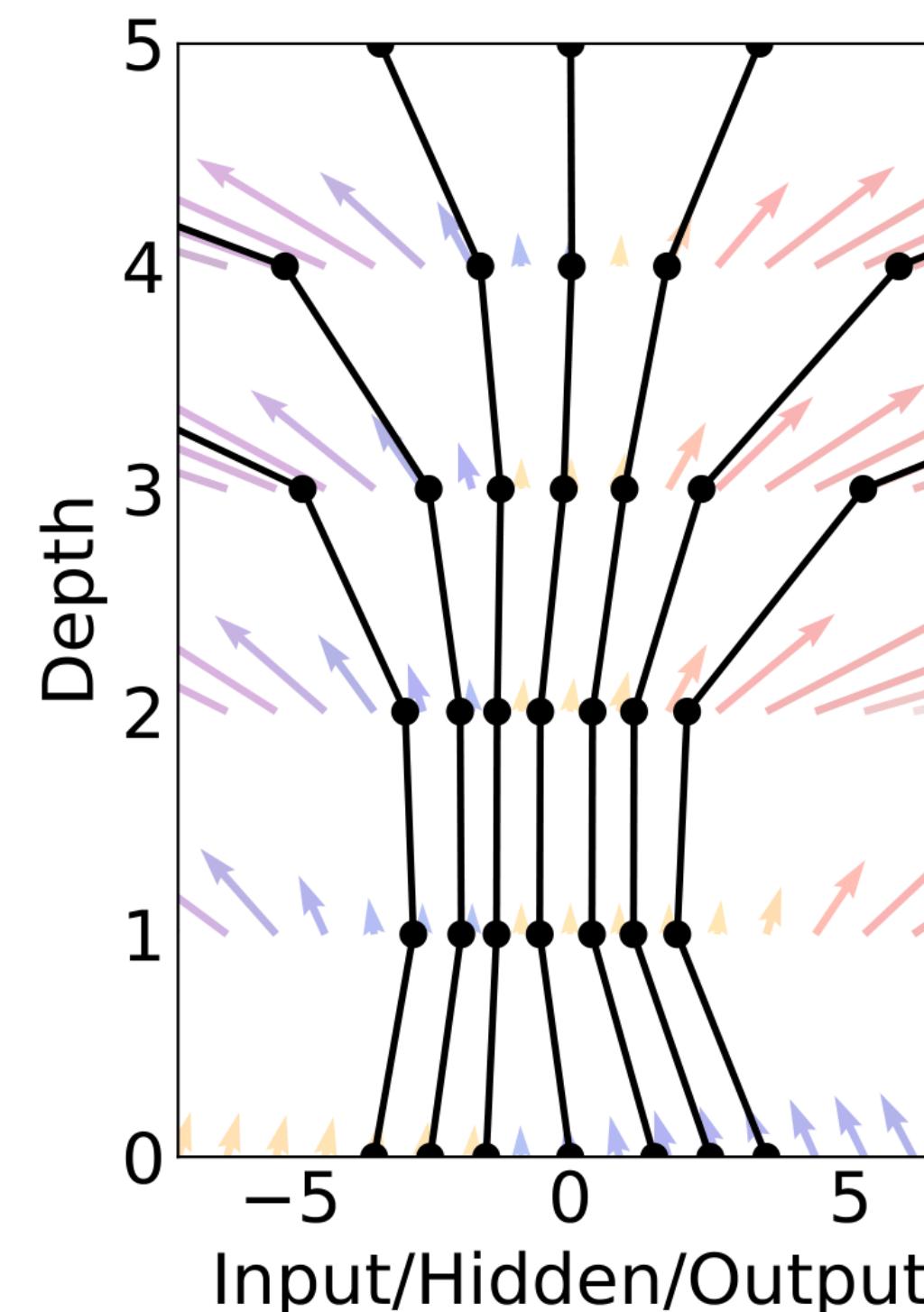


$$dx/dt = v(x)$$

Harbor et al 1705.03341
Lu et al 1710.10121,
E Commun. Math. Stat 17'

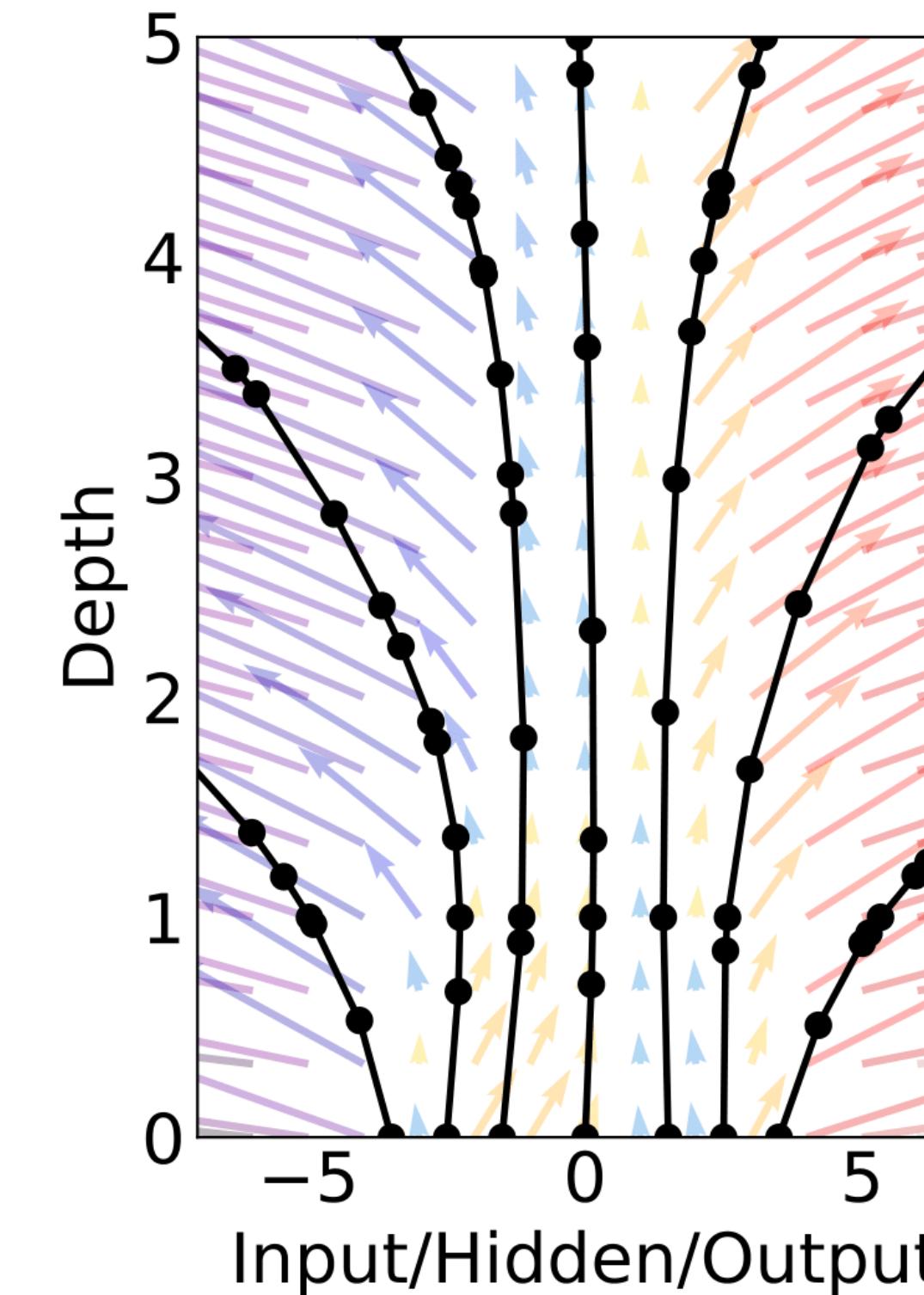
Neural Ordinary Differential Equations

Residual network



$$x_{t+1} = x_t + v(x_t)$$

ODE integration



$$dx/dt = v(x)$$

Chen et al, 1806.07366

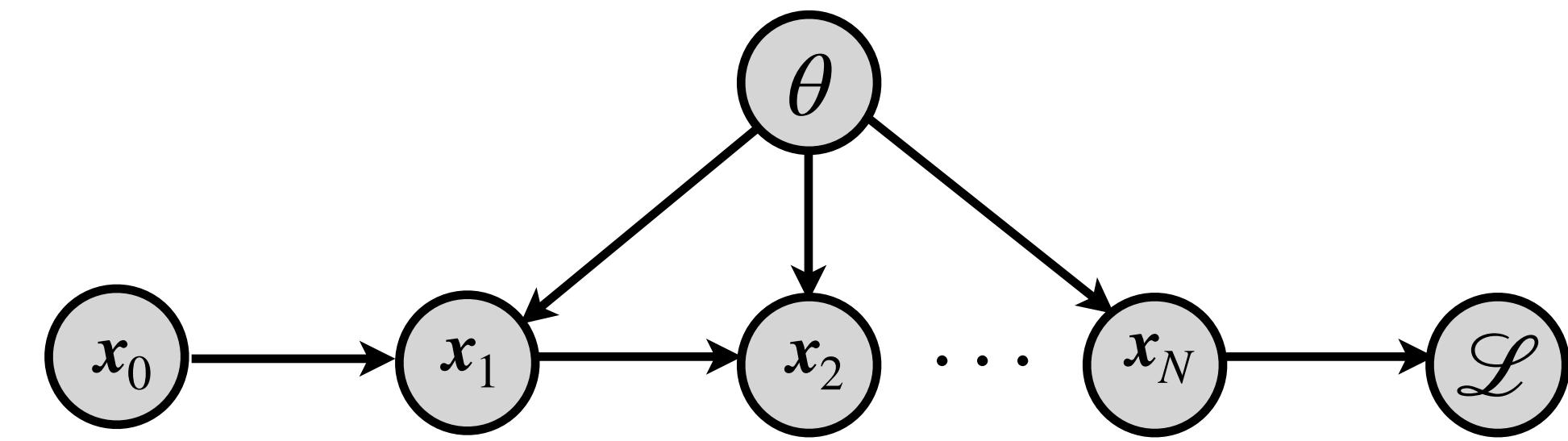
Harbor et al 1705.03341

Lu et al 1710.10121,

E Commun. Math. Stat 17'...

Backpropagate through an ODE

$$\frac{dx}{dt} = v(x, \theta, t)$$



Adjoint $\bar{x}(t) = \frac{\partial \mathcal{L}}{\partial x(t)}$ satisfies another ODE to be integrated back in time

$$\frac{d\bar{x}(t)}{dt} = -\bar{x}(t) \frac{\partial v(x, \theta, t)}{\partial x}$$

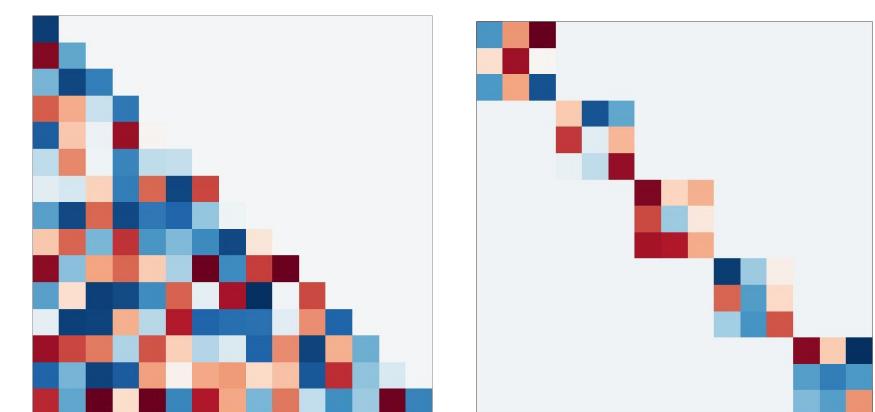
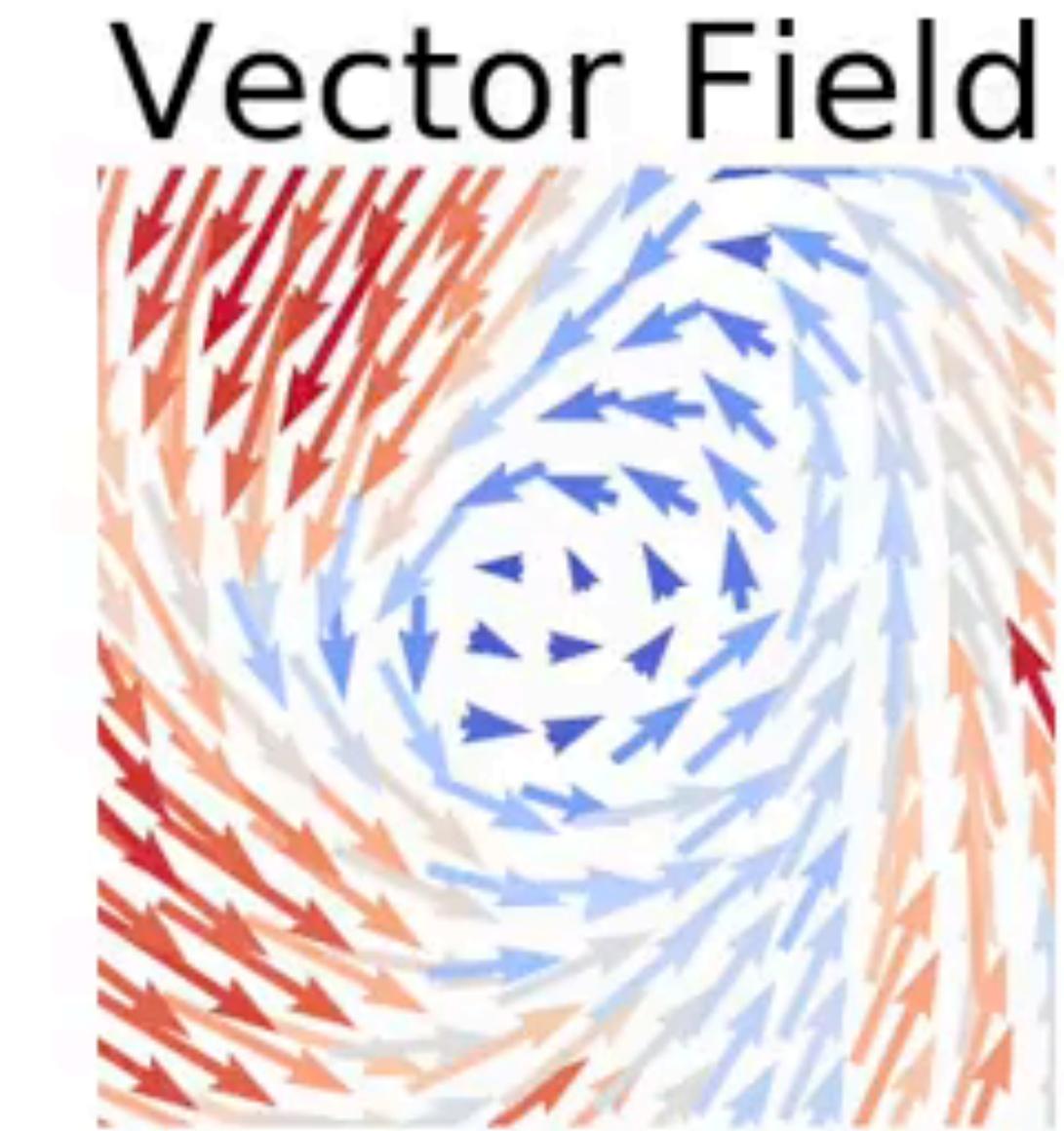
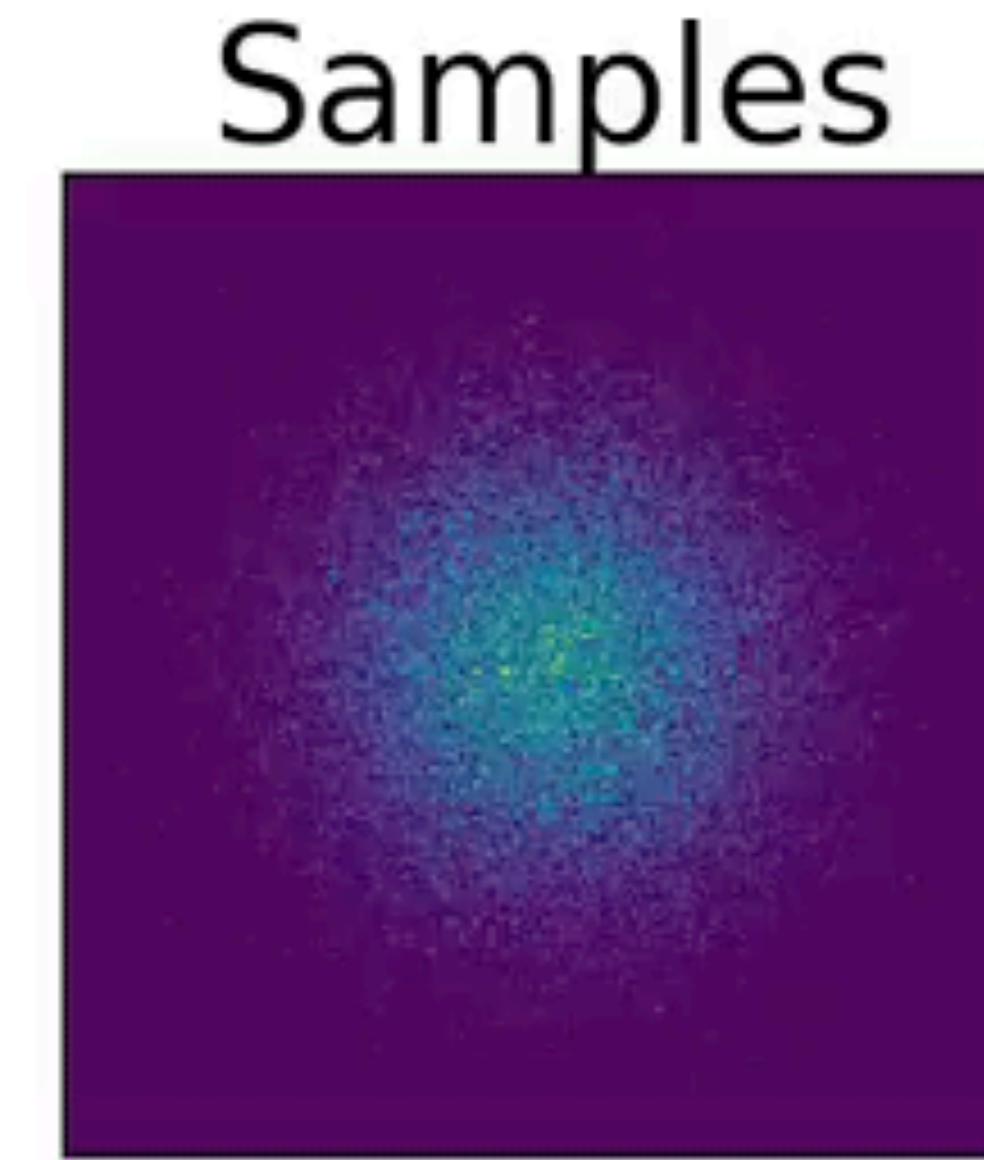
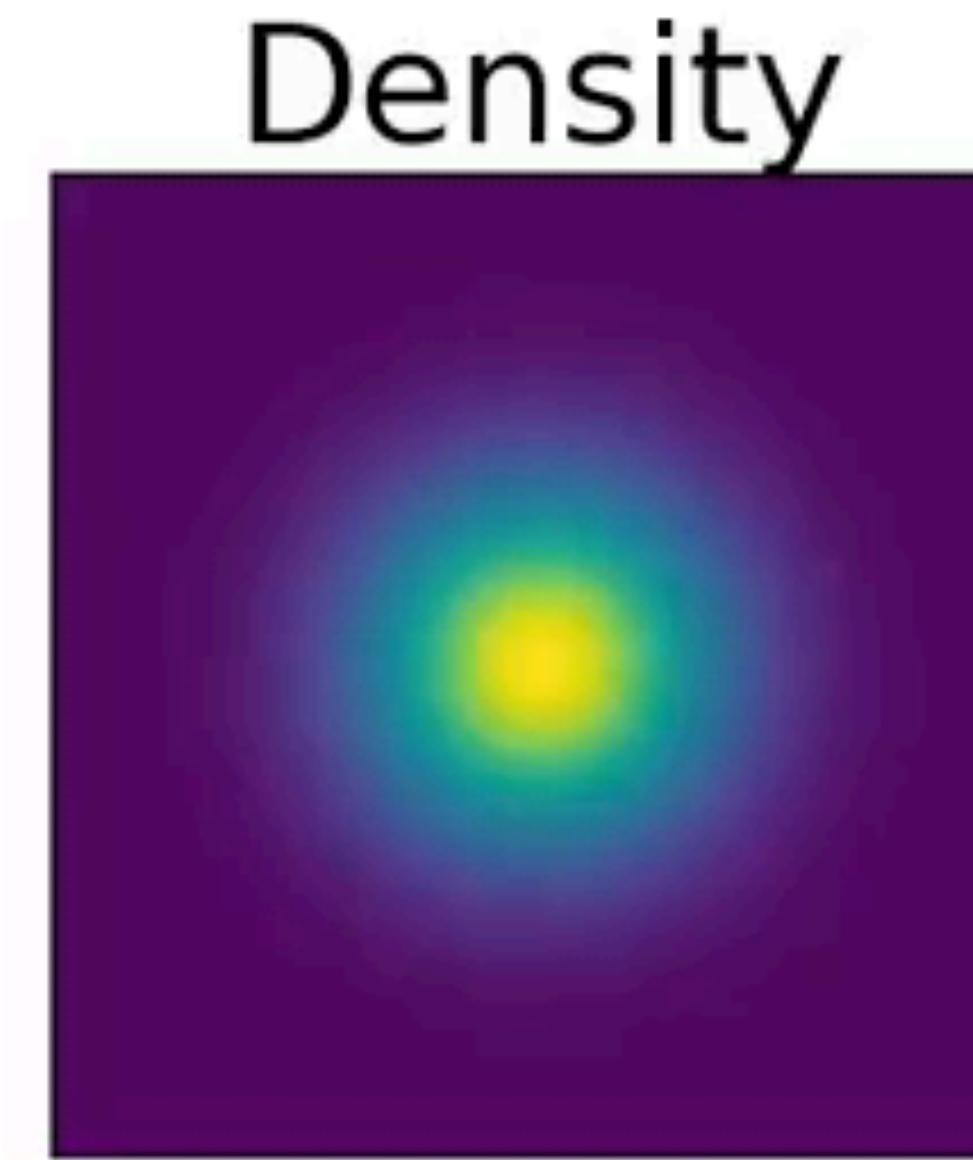
Gradient w.r.t. parameter

$$\frac{\partial \mathcal{L}}{\partial \theta} = \int_0^T dt \bar{x}(t) \frac{\partial v(x, \theta, t)}{\partial \theta}$$

Exercise:
Derive this!

Continuous normalizing flows implemented with NeuralODE

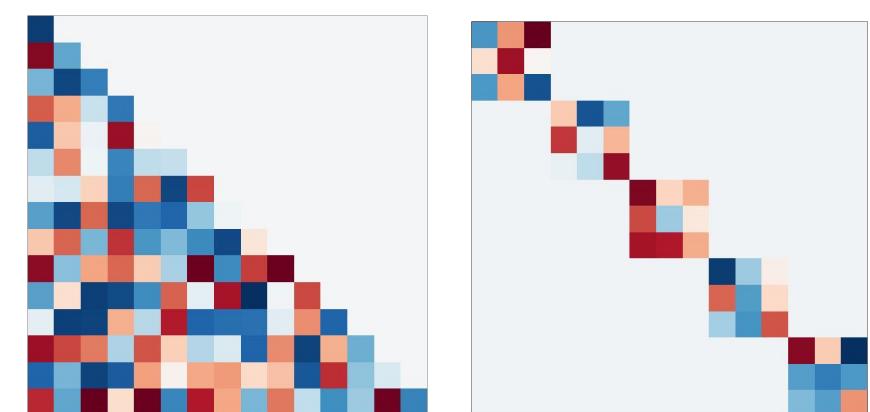
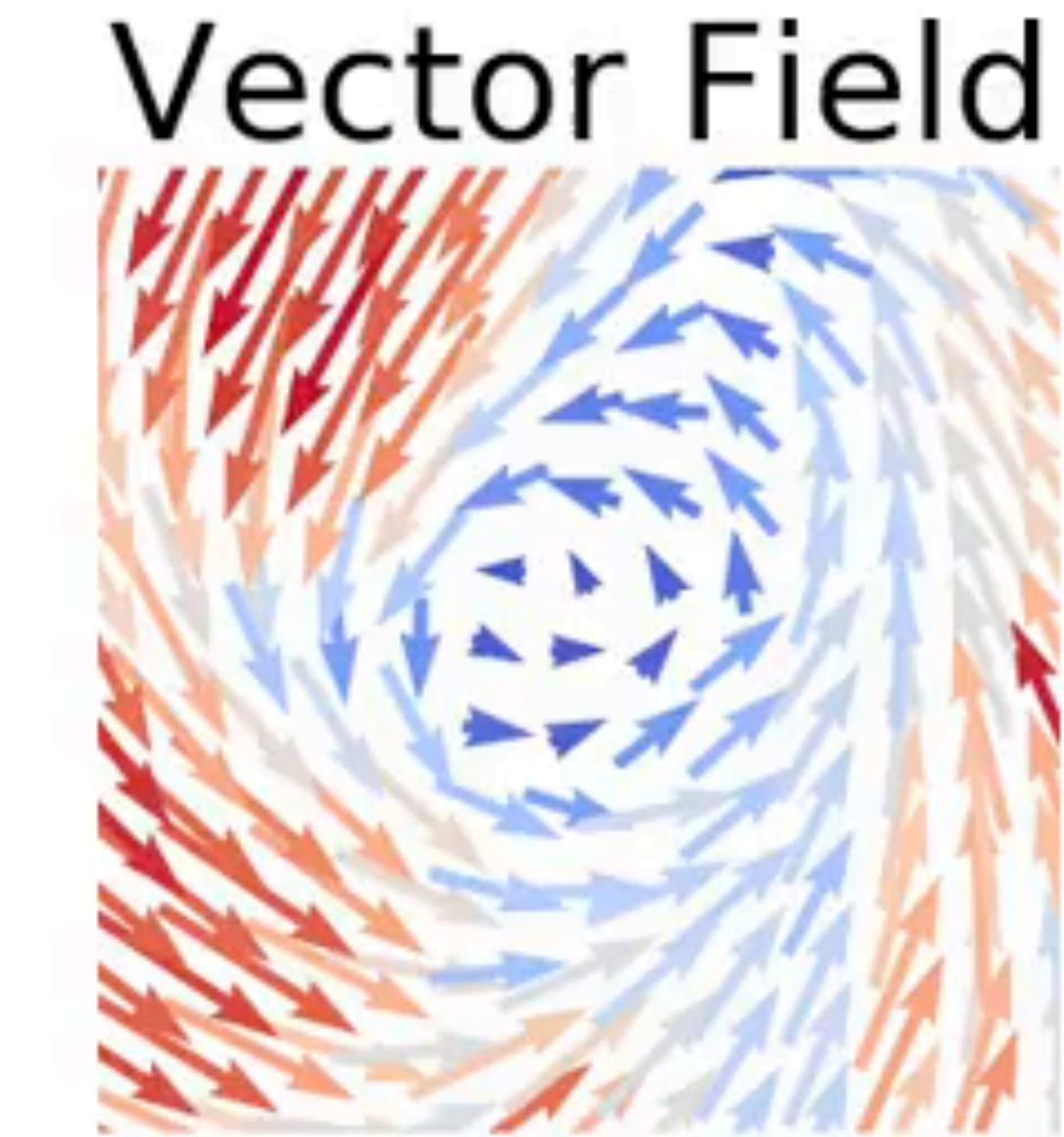
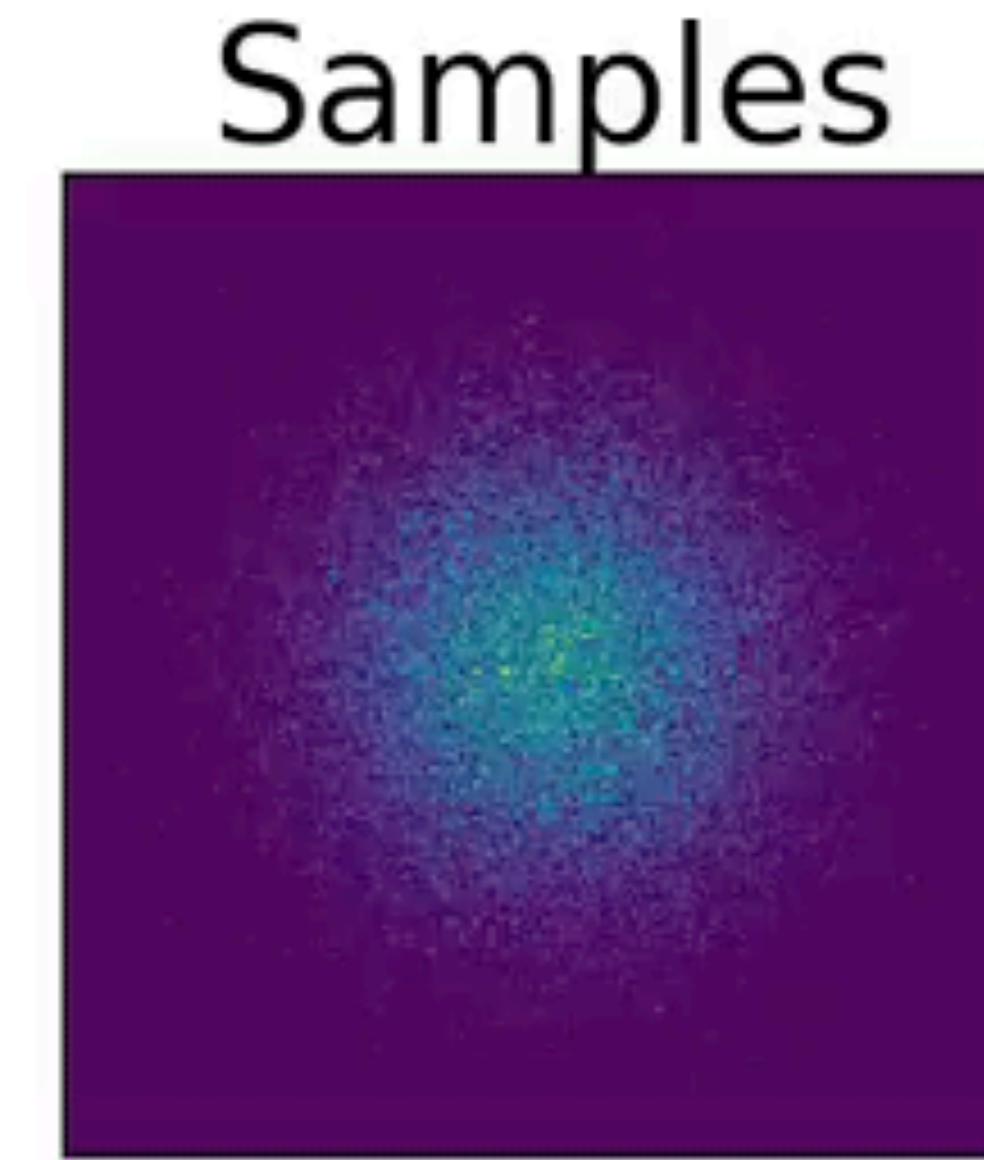
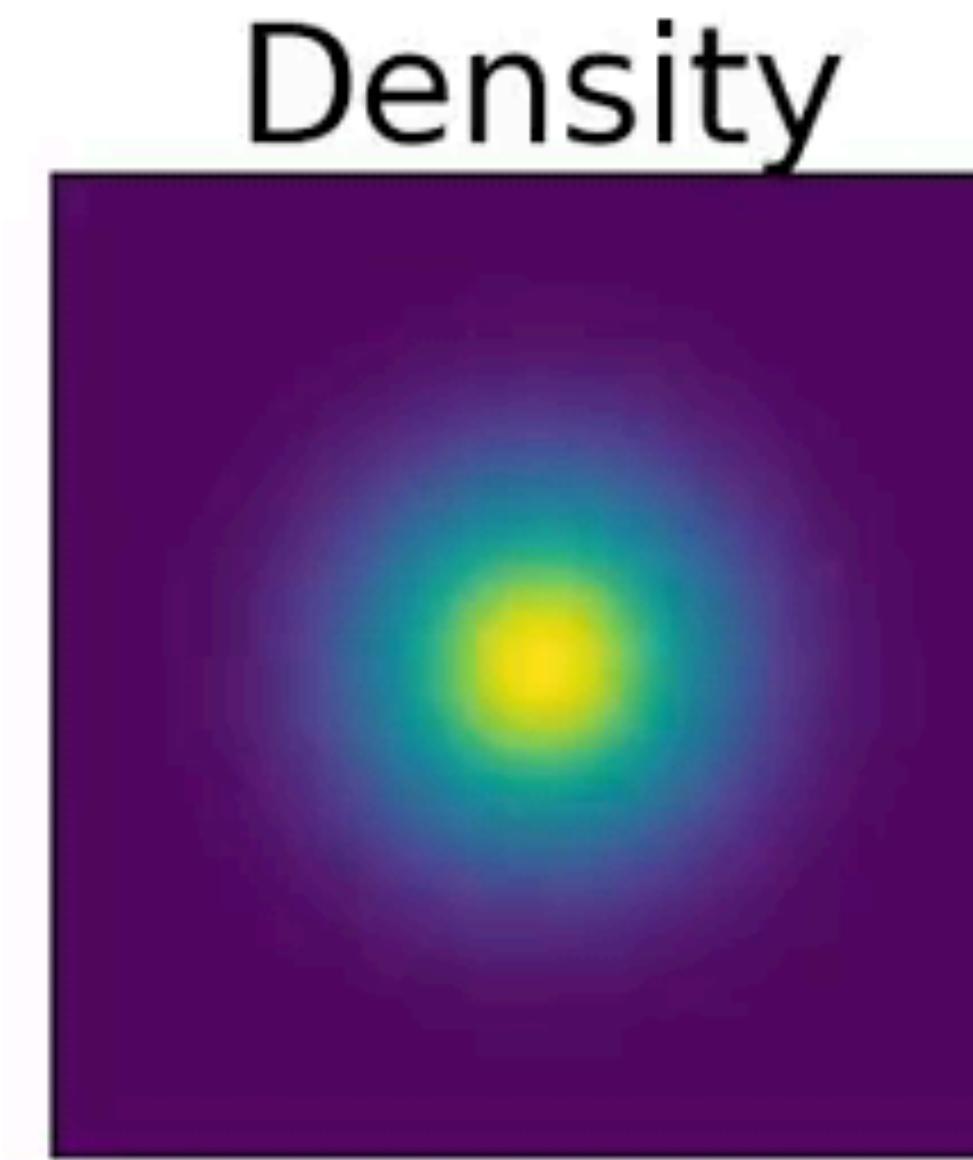
Chen et al, 1806.07366, Grathwohl et al 1810.01367



Continuous normalizing flow have no structural
constraints on the transformation Jacobian

Continuous normalizing flows implemented with NeuralODE

Chen et al, 1806.07366, Grathwohl et al 1810.01367

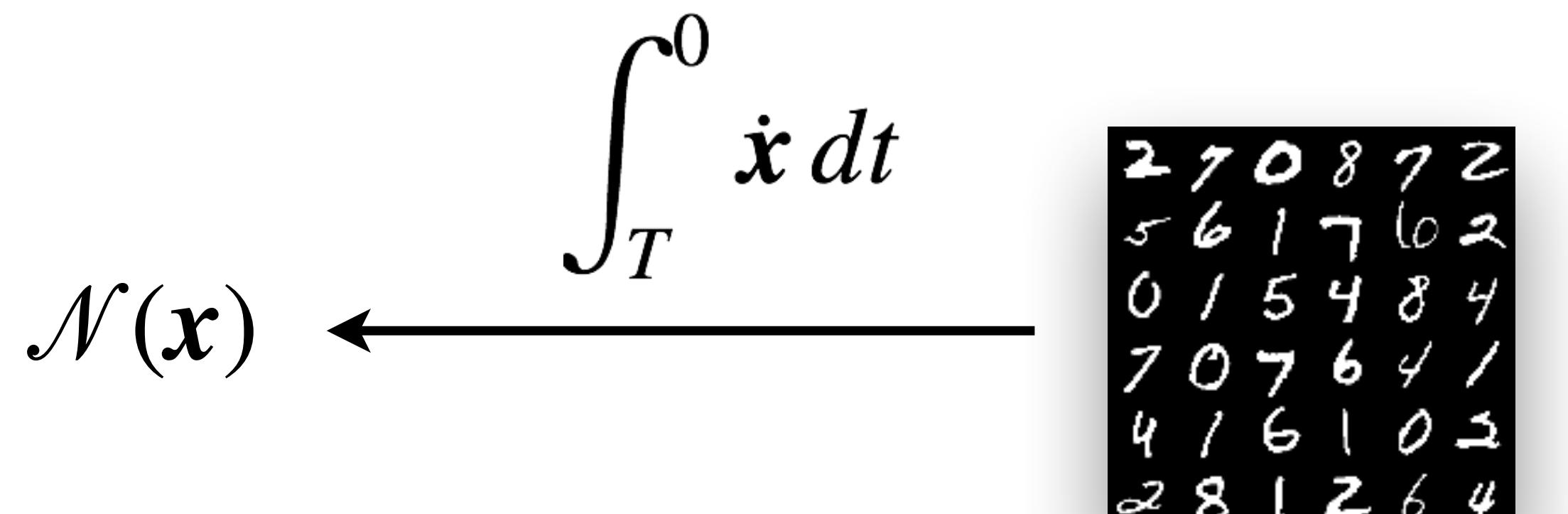


Continuous normalizing flow have no structural
constraints on the transformation Jacobian

The two use cases

Zhang, E, LW, 1809.10188

Maximum likelihood estimation



Variational free energy

$$\mathcal{N}(\mathbf{x}) \xrightarrow{\int_0^T \dot{\mathbf{x}} dt} \frac{e^{-E(\mathbf{x})}}{Z}$$

“learn from data”

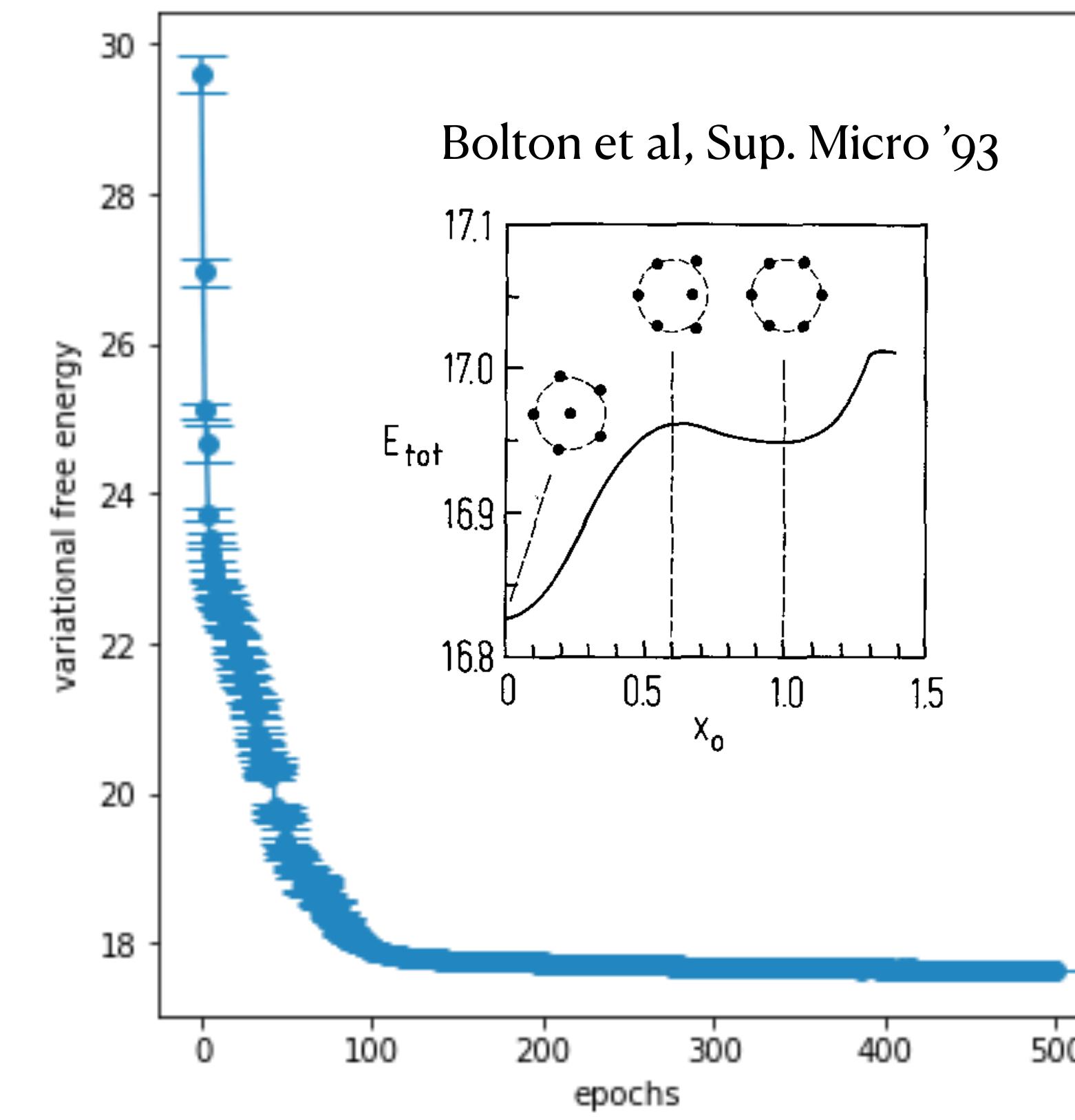
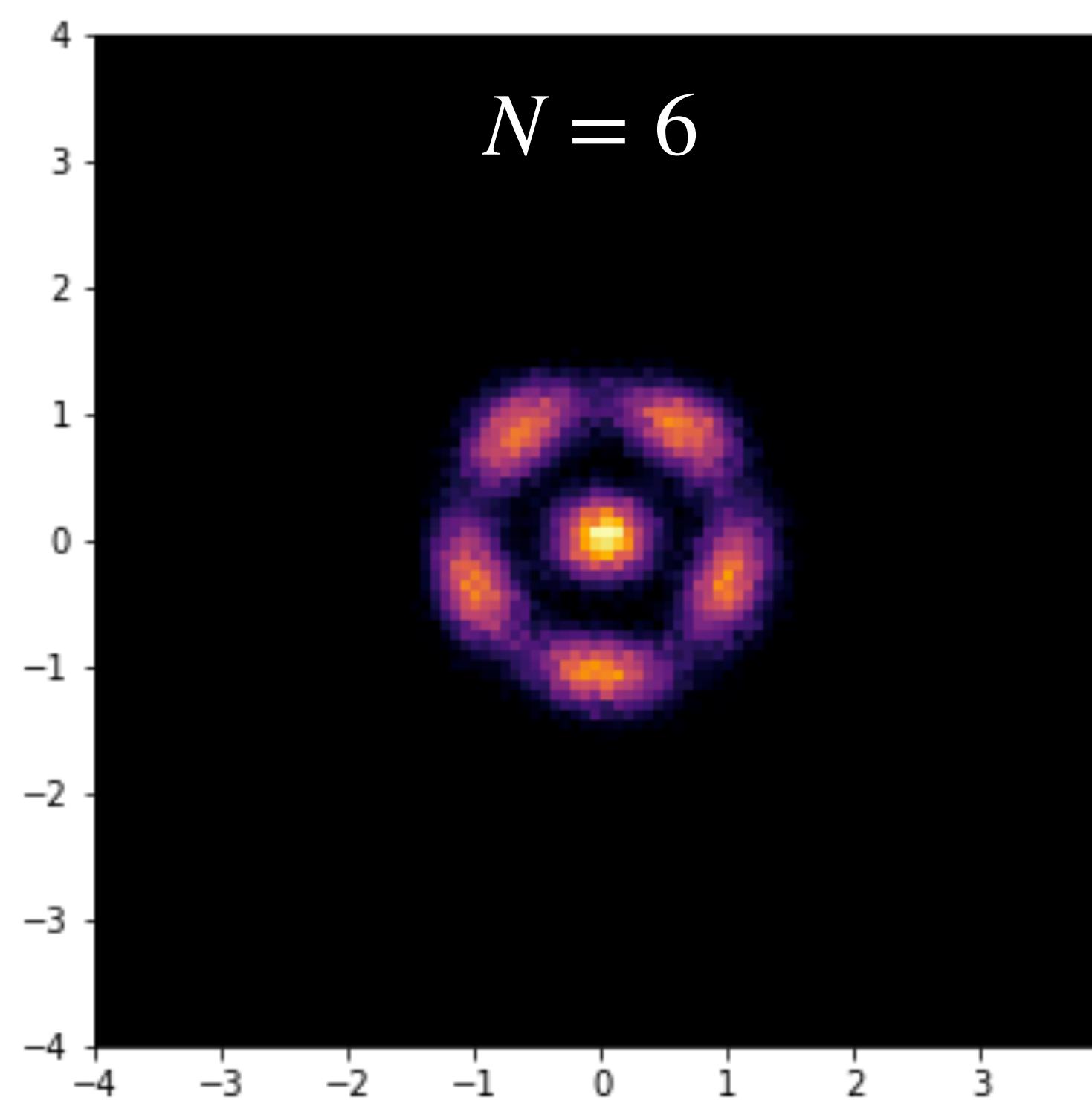
$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim \text{data}} [\ln p(\mathbf{x})]$$

“learn from Energy”

$$F = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [E(\mathbf{x}) + k_B T \ln p(\mathbf{x})]$$

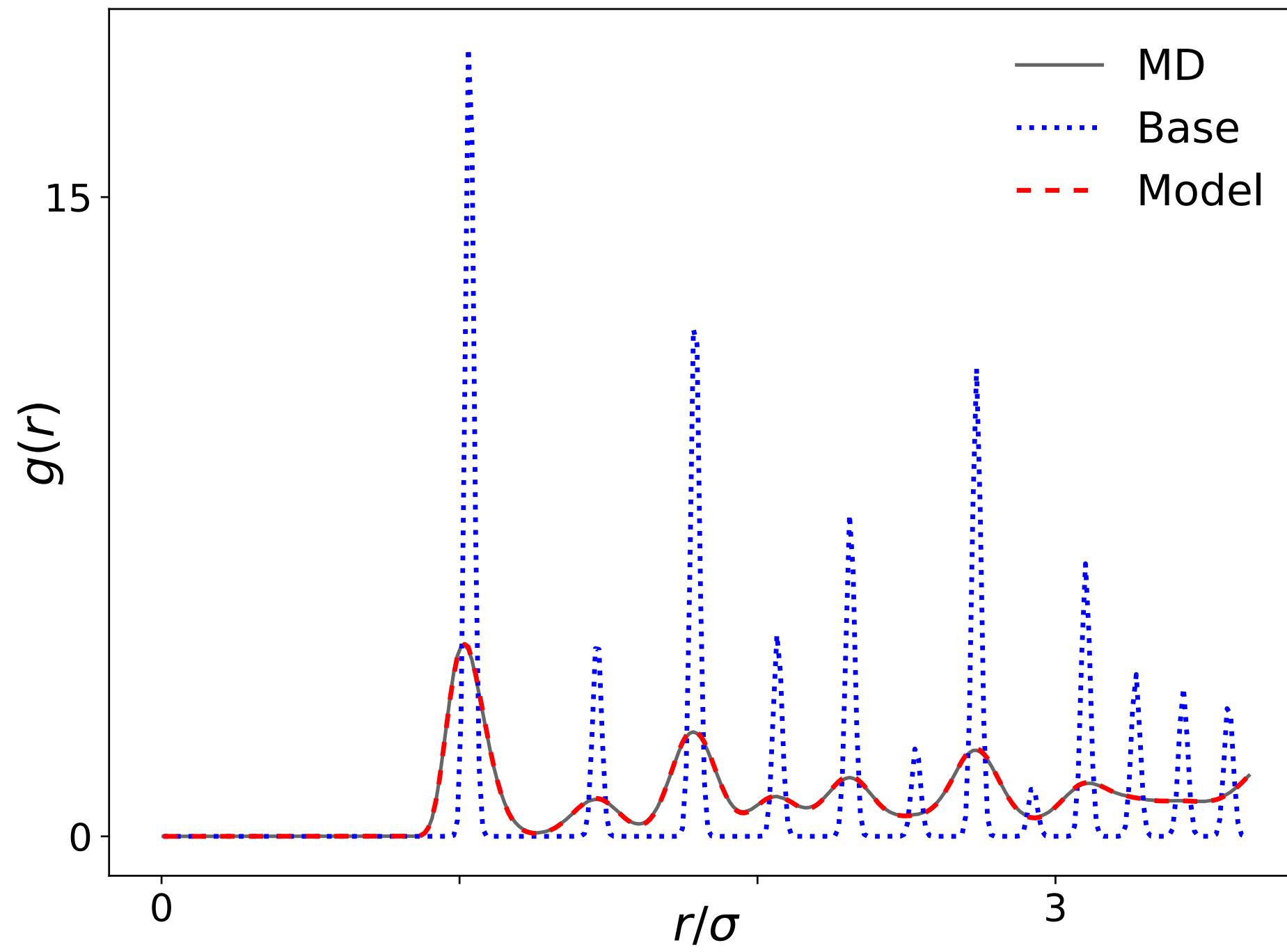
Demo: Classical Coulomb gas in a harmonic trap

$$E = \sum_{i < j} \frac{1}{|x_i - x_j|} + \sum_i x_i^2$$



Case study: Normalizing flow for atomic solids

Variational free energy with a really deep permutation equivariant flow



System	N	LFEP	LBAR	MBAR
LJ	256	3.10800(28)	3.10797(1)	3.10798(9)
LJ	500	3.12300(41)	3.12264(2)	3.12262(10)

$$\ln Z = \ln \mathbb{E}_{x \sim q(x)} [e^{-\beta E(x) - \ln q(x)}]$$

free energy perturbation (Zwanzig 1954)

$$\ln Z_B - \ln Z_A = \ln \mathbb{E}_A [e^{-\beta(E_B - E_A)}]$$

Normalizing flow for atomic solids

F. Hardware details and computational cost

For our flow experiments, we used 16 A100 GPUs to train each model on the bigger systems (512-particle mW and 500-particle LJ). It took approximately 3 weeks of training to reach convergence of the free-energy estimates. Obtaining 2M samples for evaluation took approximately 12 hours on 8 V100 GPUs for each of these models.

For each baseline MBAR estimate, we performed 100 separate simulations for LJ and 200 for mW, corresponding to the number of stages employed. These simulations were performed with LAMMPS [8] and each of them ran on multiple CPU cores communicating via MPI. We used 4 cores for the 64-particle and 216-particle mW experiments and 8 cores for all other systems. The MD simulations completed after approximately 11 and 14 hours for LJ (256 and 500 particles), and 7, 20 and 48 hours for mW (64, 216 and 512 particles). To evaluate the energy matrix for a single MBAR

Heavy lifting is mostly due to back-and-forth simulation of deep equivariant flow

Training: Monte Carlo Gradient Estimators

Review: 1906.10652

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)]$$

Reinforcement learning

Variational inference

Variational Monte Carlo

Variational quantum algorithms

...

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)] = \mathbb{E}_{x \sim p_{\theta}} [f(x) \nabla_{\theta} \ln p_{\theta}(x)]$$

Score function estimator (REINFORCE)

Pathwise estimator (Reparametrization trick) $x = g_{\theta}(z)$

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)] = \mathbb{E}_{z \sim \mathcal{N}(z)} [\nabla_{\theta} f(g_{\theta}(z))]$$

10.1 Guidance in Choosing Gradient Estimators

With so many competing approaches, we offer our rules of thumb in choosing an estimator, which follow the intuition we developed throughout the paper:

- If our estimation problem involves continuous functions and measures that are continuous in the domain, then using the pathwise estimator is a good default. It is relatively easy to implement and a default implementation, one without other variance reduction, will typically have variance that is low enough so as not to interfere with the optimisation.
- If the cost function is not differentiable or a black-box function then the score-function or the measure-valued gradients are available. If the number of parameters is low, then the measure-valued gradient will typically have lower variance and would be preferred. But if we have a high-dimensional parameter set, then the score function estimator should be used.
- If we have no control over the number of times we can evaluate a black-box cost function, effectively only allowing a single evaluation of it, then the score function is the only estimator of the three we reviewed that is applicable.
- The score function estimator should, by default, always be implemented with at least a basic variance reduction. The simplest option is to use a baseline control variate estimated with a running average of the cost value.
- When using the score-function estimator, some attention should be paid to the dynamic range of the cost function and its variance, and to find ways to keep its value bounded within a reasonable range, e.g., transforming the cost so that it is zero mean, or using a baseline.
- For all estimators, track the variance of the gradients if possible and address high variance by using a larger number of samples from the measure, decreasing the learning rate, or clipping the gradient values. It may also be useful to restrict the range of some parameters to avoid extreme values, e.g., by clipping them to a desired interval.
- The measure-valued gradient should be used with some coupling method for variance reduction. Coupling strategies that exploit relationships between the positive and negative components of the density decomposition, and which have shared sampling paths, are known for the commonly-used distributions.
- If we have several unbiased gradient estimators, a convex combination of them might have lower variance than any of the individual estimators.
- If the measure is discrete on its domain then the score-function or measure-valued gradient are available. The choice will again depend on the dimensionality of the parameter space.
- In all cases, we strongly recommend having a broad set of tests to verify the unbiasedness of the gradient estimator when implemented.

Mohamed et al, 1906.10652

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)]$$

When to use which ?

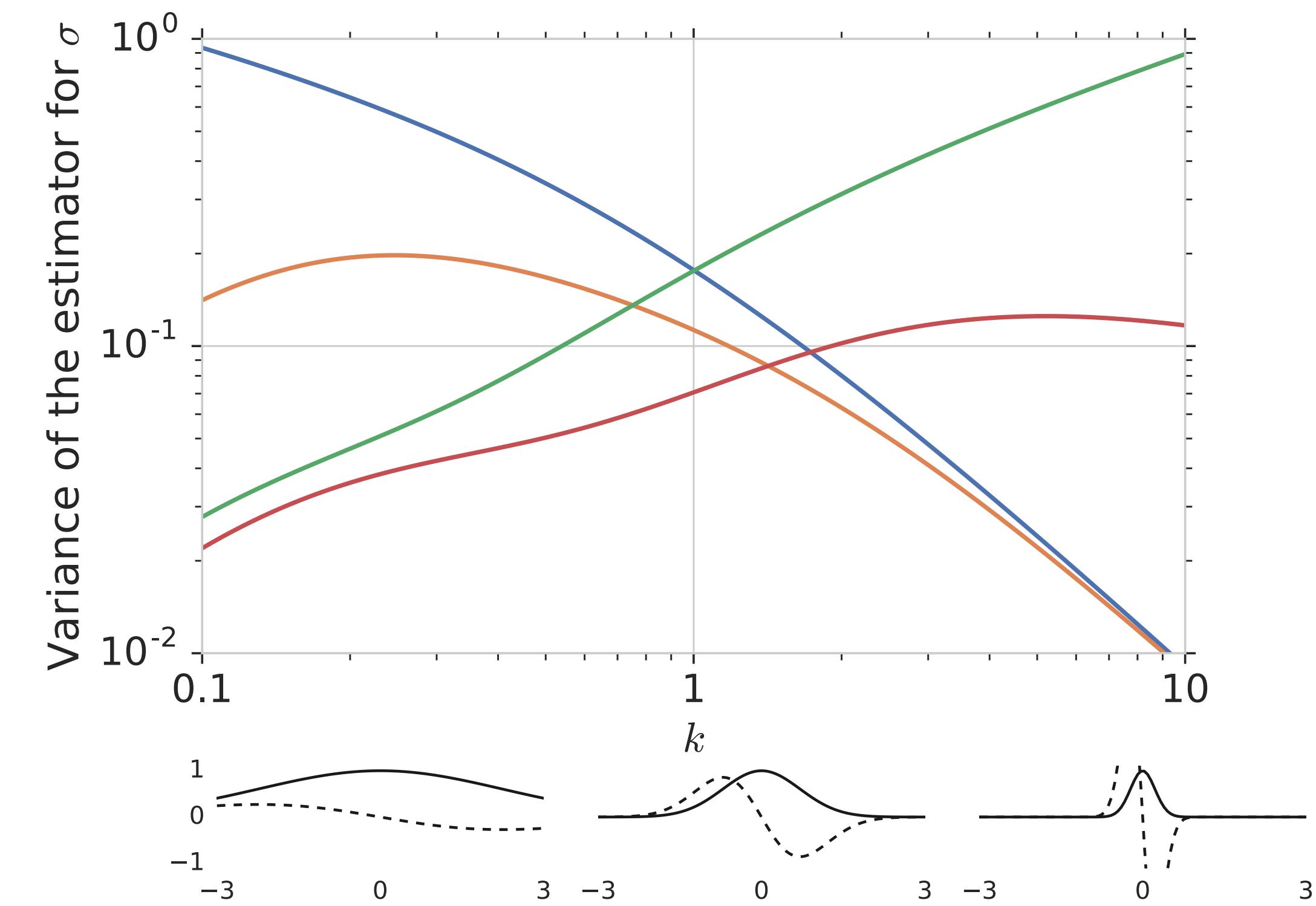
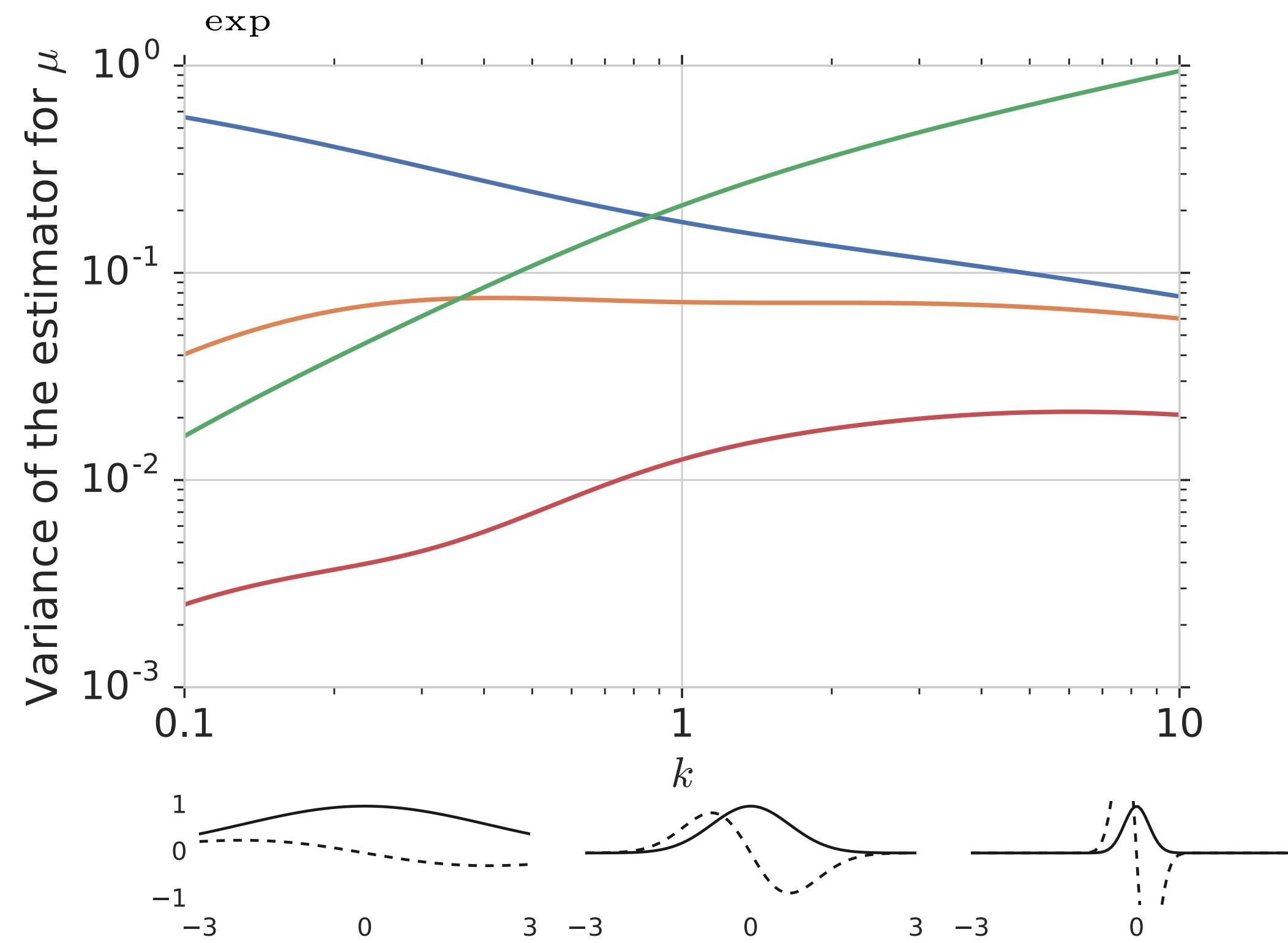
More discussions

Roeder et al, 1703.09194

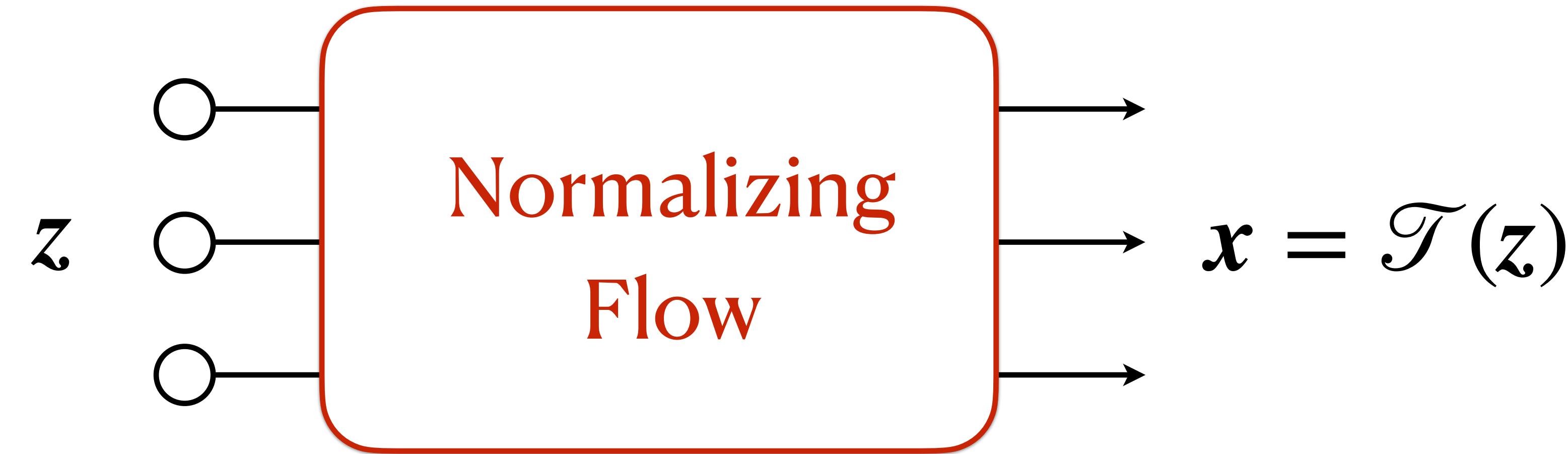
Vaitl et al 2206.09016, 2207.08219

$$\eta = \nabla_{\theta} \int \mathcal{N}(x|\mu, \sigma^2) f(x; k) dx; \quad \theta \in \{\mu, \sigma\}$$

— Score function — Score function + variance reduction — Pathwise — Measure-valued + variance reduction
— Value of the cost --- Derivative of the cost



Symmetries



Invariance

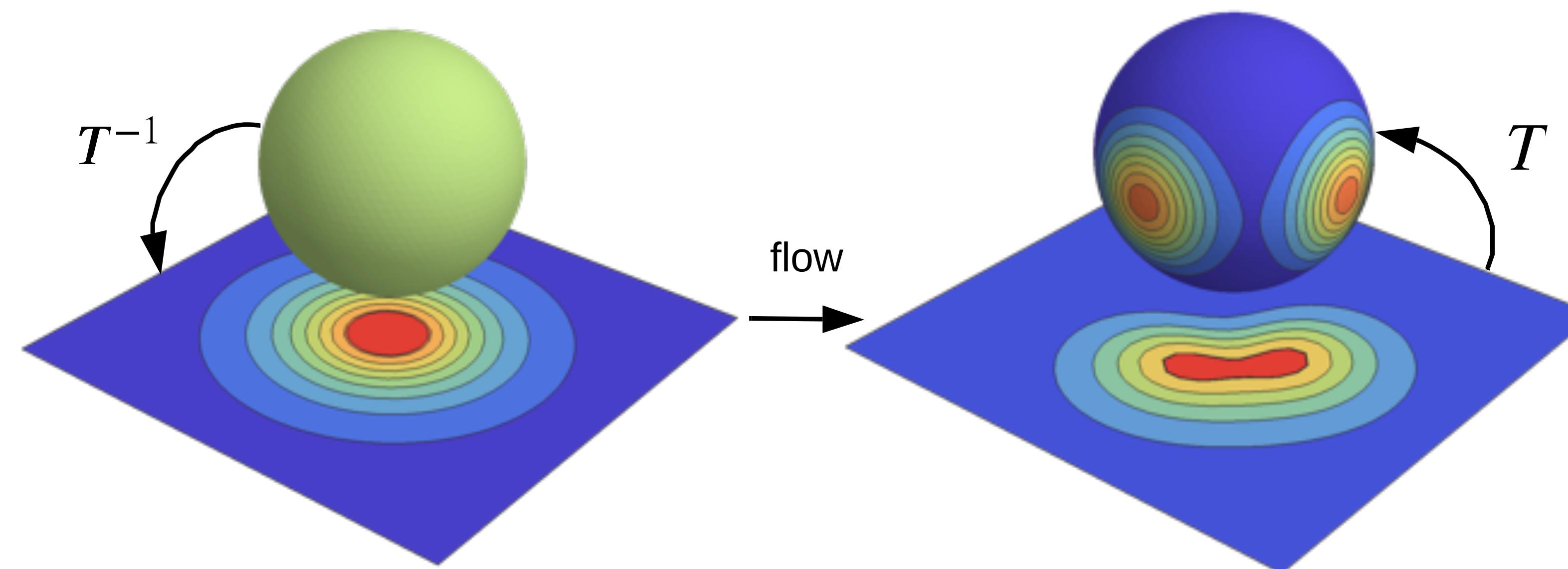
$$\rho(gx) = \rho(x)$$

Equivariance

$$\mathcal{T}(gz) = g\mathcal{T}(z)$$

Spatial symmetries, permutation symmetries, gauge symmetries...

Flow on manifolds

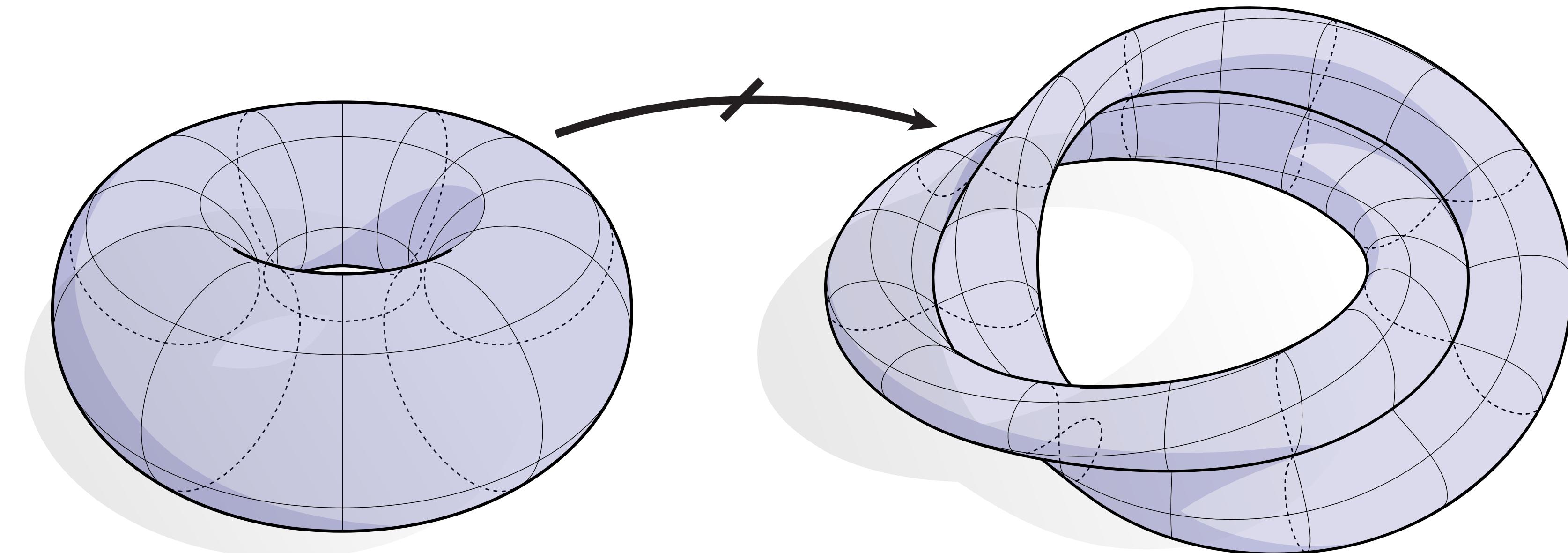


Periodic variables, gauge fields, ...

Gemici et al 1611.02304, Rezende et al, 2002.02428, Boyda et al, 2008.05456

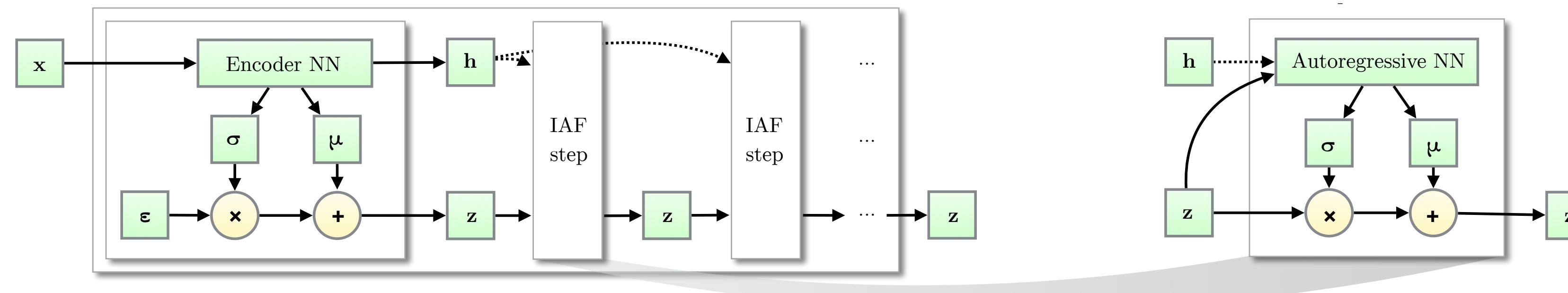
Neural ODE on manifolds, Falorsi et al, 2006.06663, Lou et al, 2006.10254, Mathieu et al, 2006.10605

Obstructions

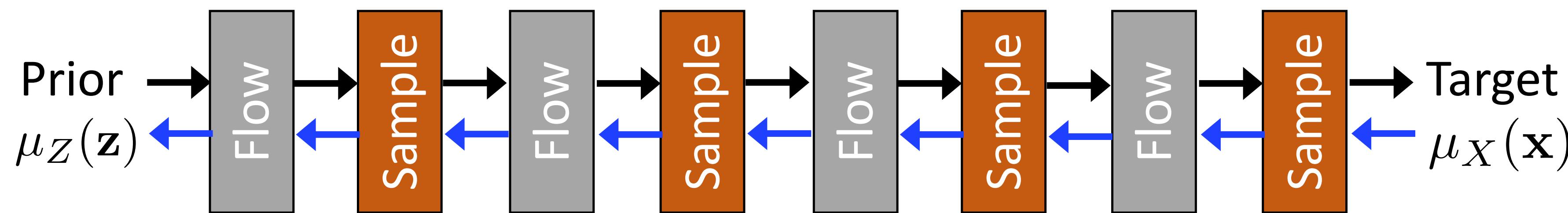


Dupont et al 1904.01681, Cornish et al, 1909.13833, Zhang et al, 1907.12998, Zhong et al, 2006.00392...

Mix with other approaches



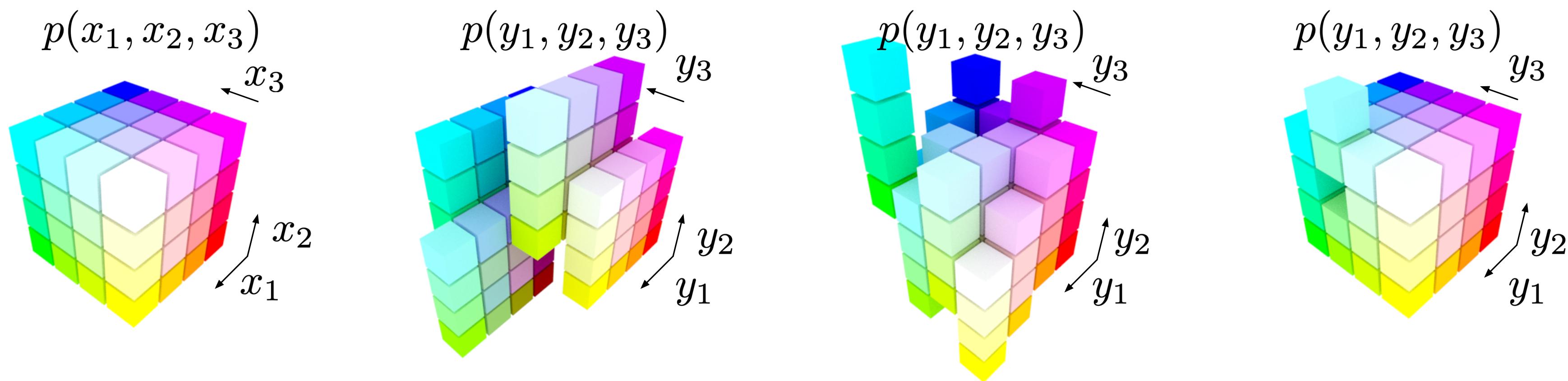
Kingma et al, 1606.04934, ...



Levy et al, 1711.09268, Wu et al 2002.06707, ...

Discrete flows

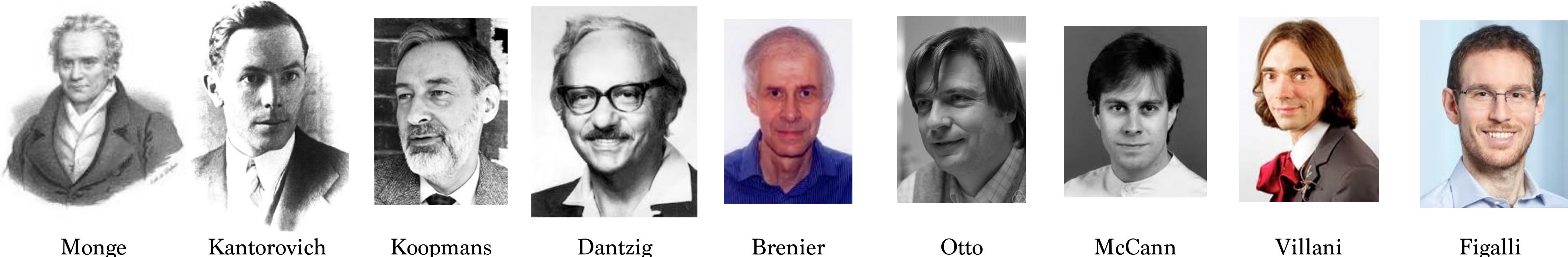
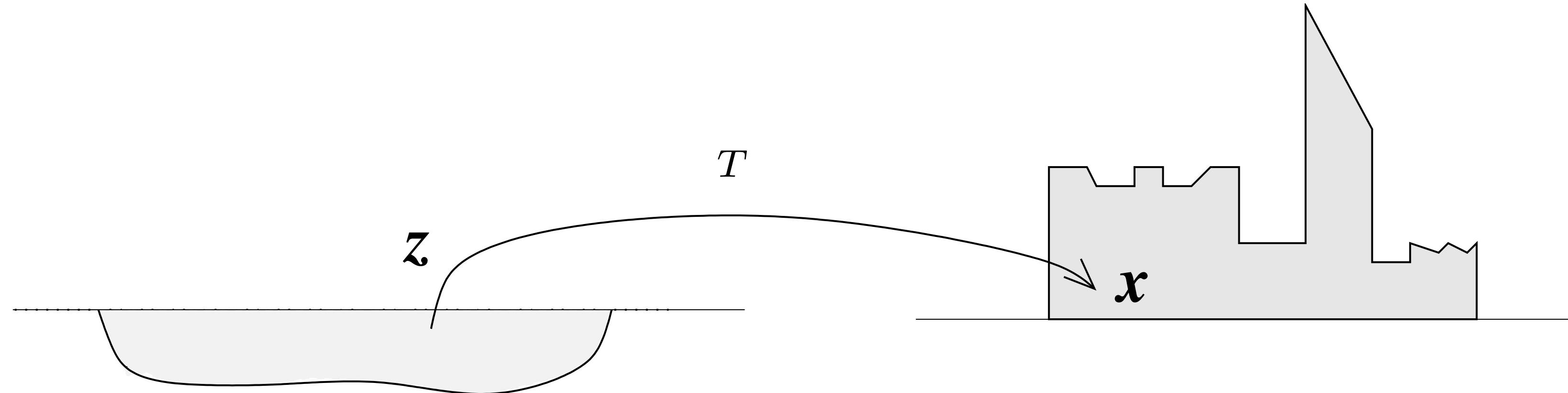
$$p(\mathbf{x}) = p(\mathbf{y} = \mathcal{T}(\mathbf{x}))$$



Tran et al, 1905.10347, Hoogeboom et al, 1905.07376, van den Berg 2006.12459

Optimal Transport Theory

Monge problem (1781): How to transport earth with optimal cost ?



Nobel Prize in Economics '75

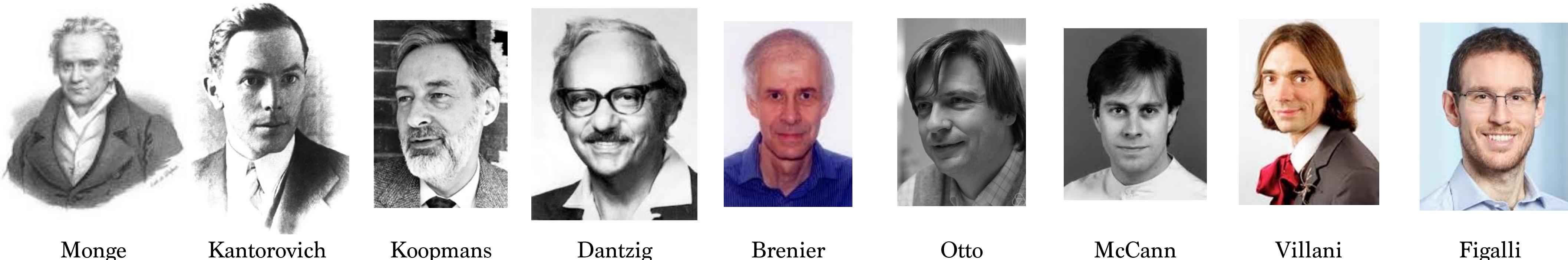
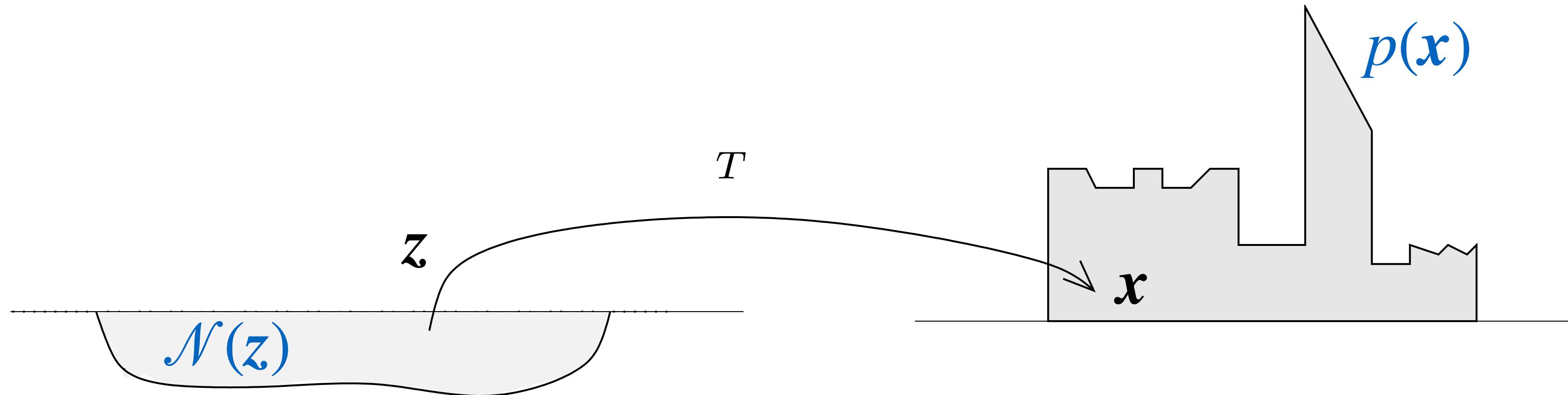
Fields Medal '10

Fields Medal '18

from Cuturi, Solomon NISP 2017 tutorial

Optimal Transport Theory

Monge problem (1781): How to transport earth with optimal cost ?



Monge

Kantorovich

Koopmans

Dantzig

Brenier

Otto

McCann

Villani

Figalli

Nobel Prize in Economics '75

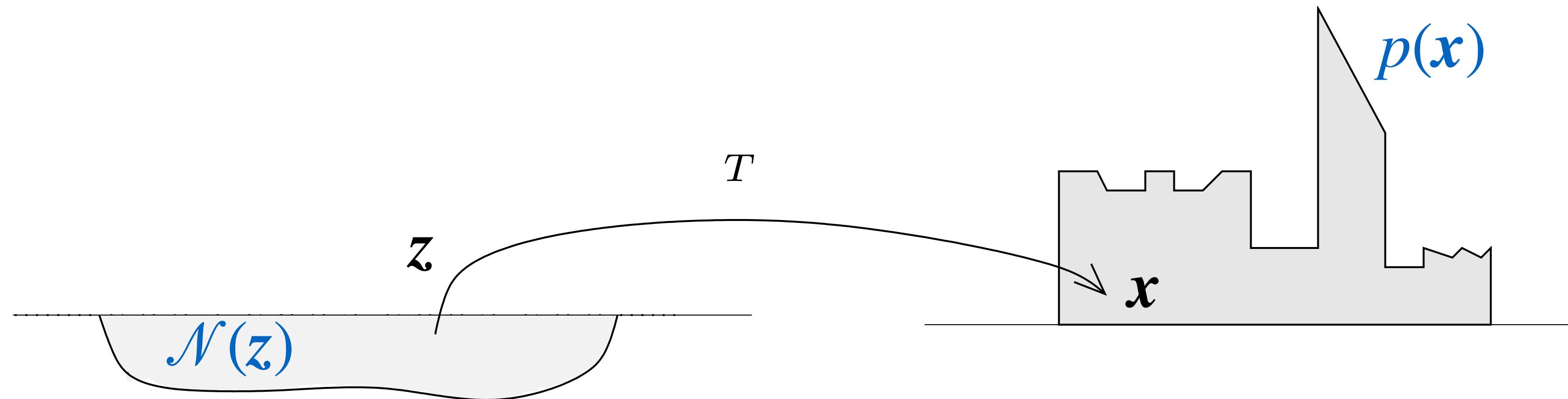
Fields Medal '10

Fields Medal '18

from Cuturi, Solomon NISP 2017 tutorial

Optimal Transport Theory

Monge problem (1781): How to transport earth with optimal cost ?



Brenier theorem (1991)

Under certain conditions
the optimal map is

$$z \mapsto x = \nabla u(z)$$

Monge-Ampère Equation

$$\frac{\mathcal{N}(z)}{p(\nabla u(z))} = \det \left(\frac{\partial^2 u}{\partial z_i \partial z_j} \right)$$

Monge-Ampère Flow

Zhang, E, LW 1809.10188

[wangleiphy/MongeAmpereFlow](#)

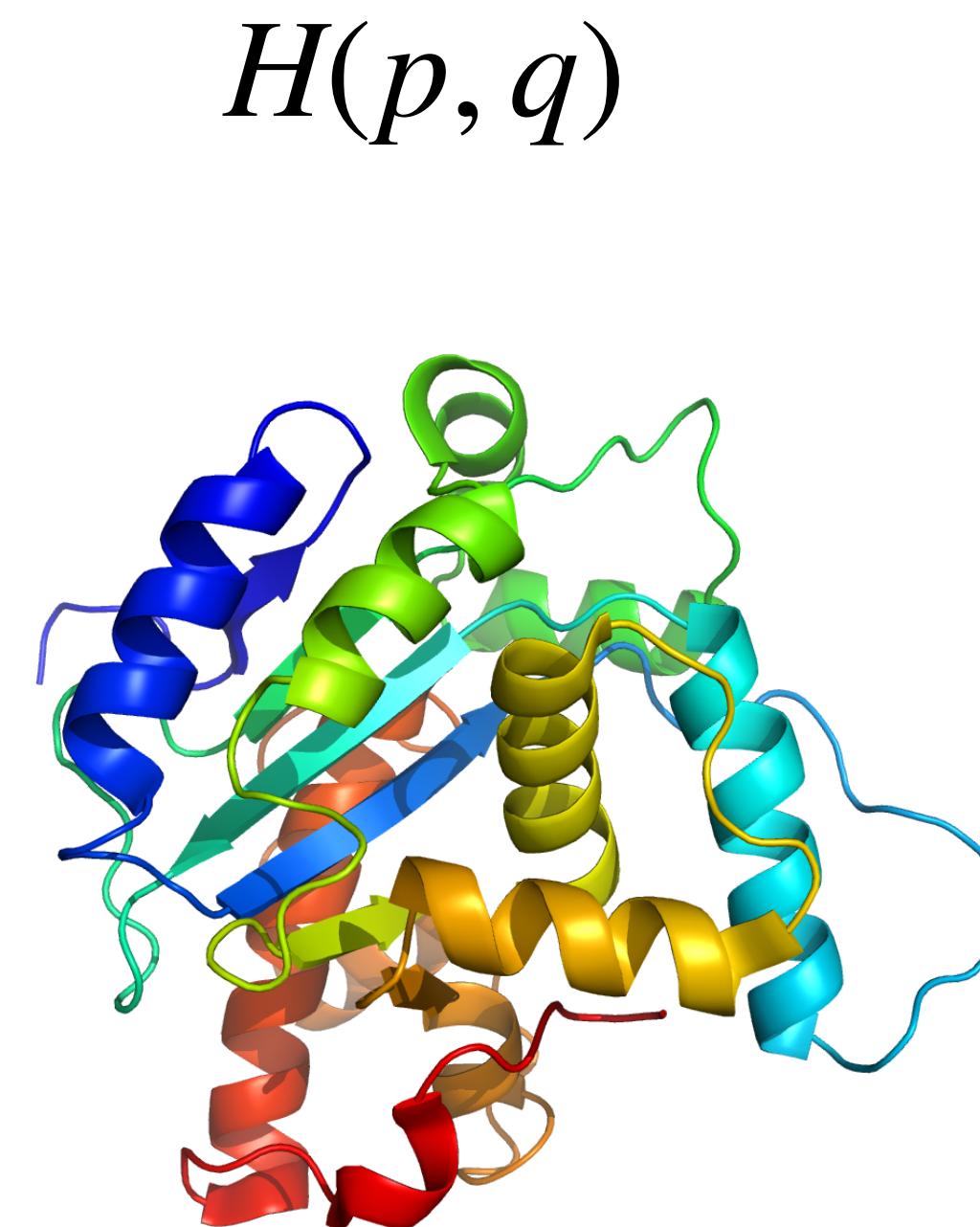
$$\frac{\partial \rho(\mathbf{x}, t)}{\partial t} + \nabla \cdot [\rho(\mathbf{x}, t) \nabla \varphi] = 0$$

- ① Drive the flow with an “irrotational” velocity field
- ② Impose symmetry to the scalar valued potential for symmetric generative model

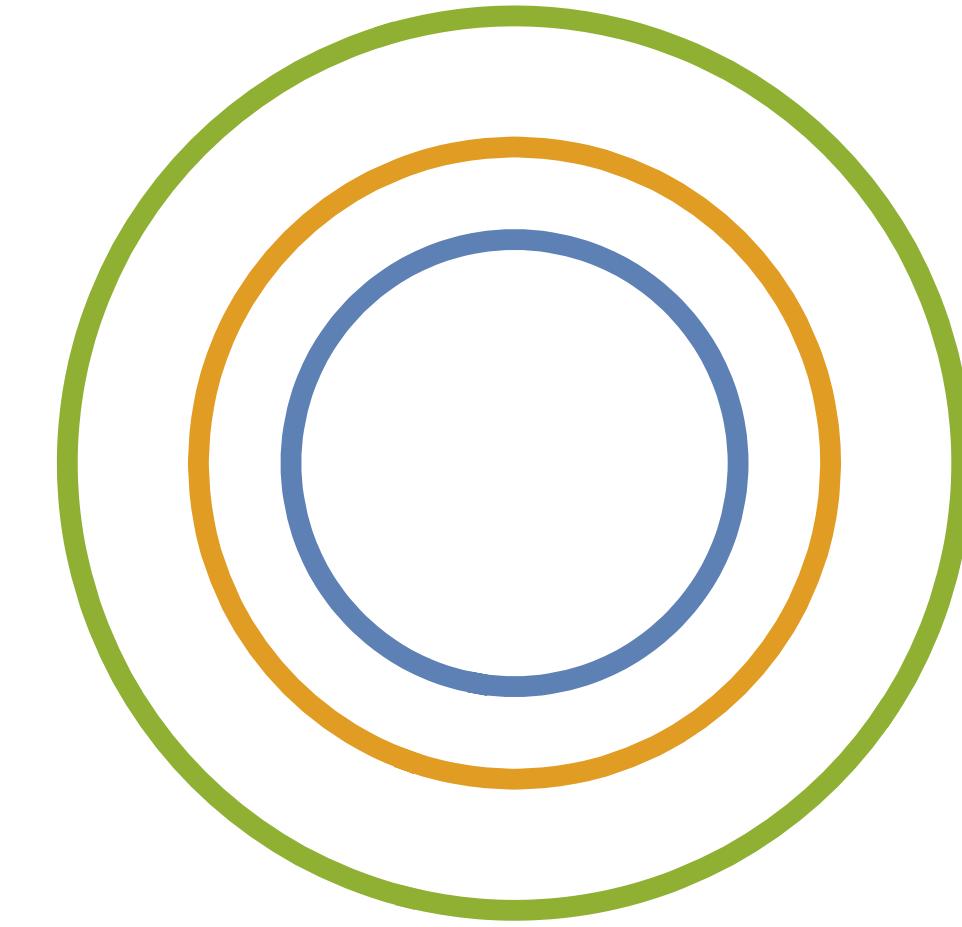
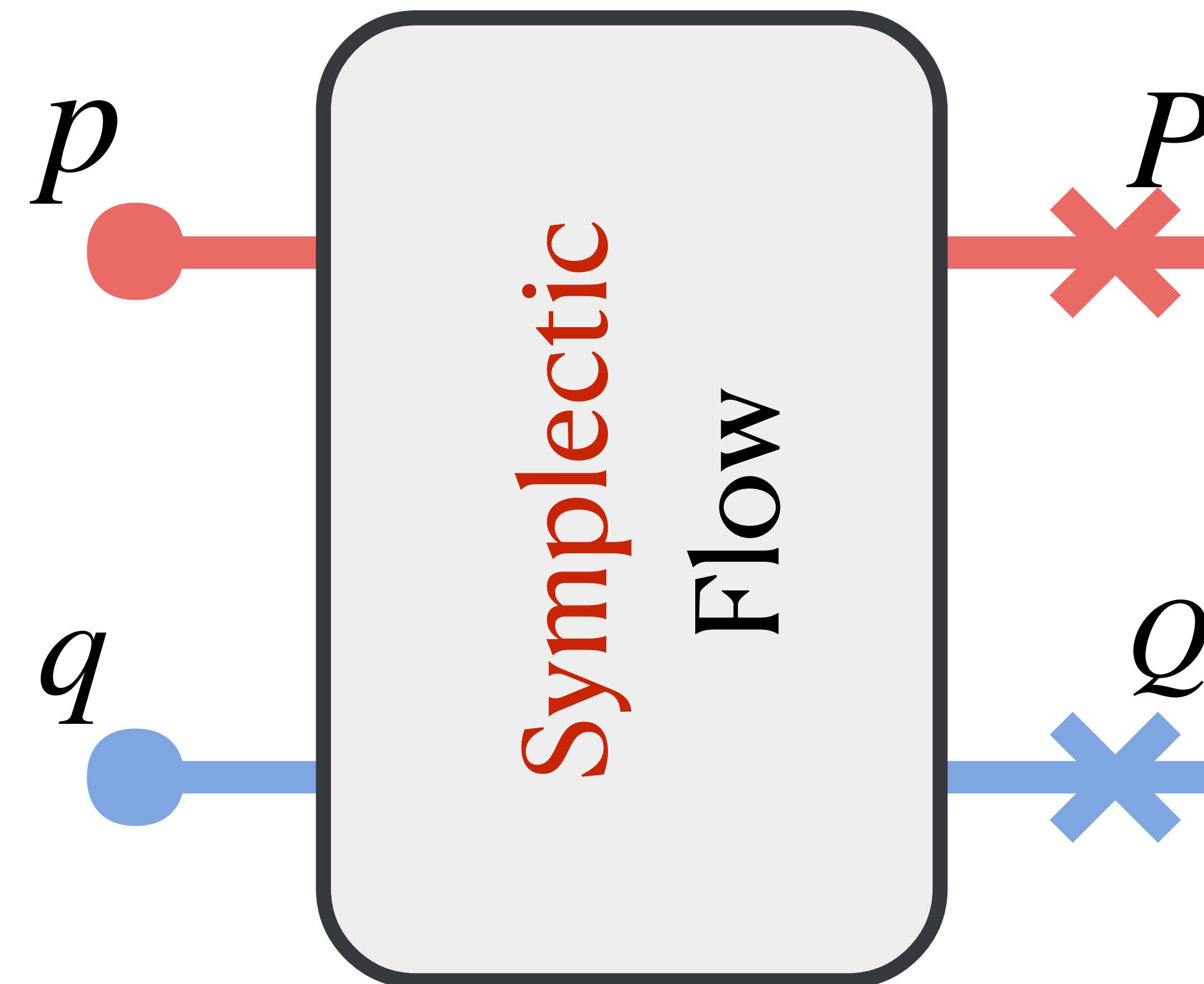
$$\varphi(g\mathbf{x}) = \varphi(\mathbf{x}) \implies \rho(g\mathbf{x}) = \rho(\mathbf{x})$$

Neural Canonical Transformations

Li, Dong, Zhang, LW, PRX '20 [lio12589/neuralCT](#)



physical space

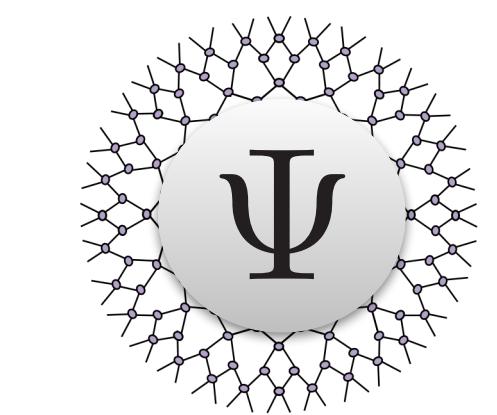


$$K(P, Q) = \sum_k \frac{P_k^2 + \omega_k^2 Q_k^2}{2}$$

Learn harmonic frequencies of the base to identify slow collective modes

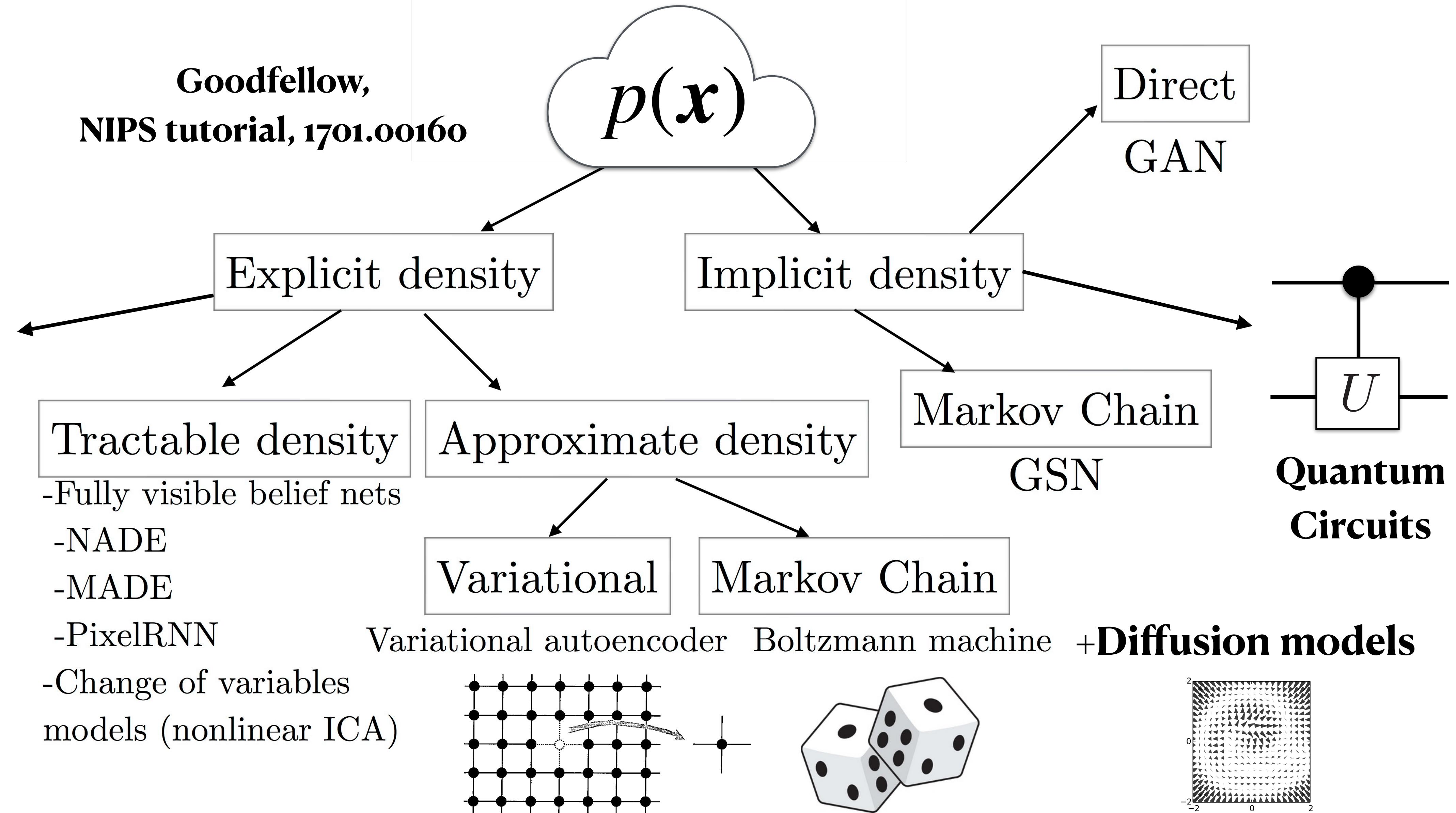
See Bondesan et al 1906.04645, Ishikawa et al 2103.00372 for investigations on integrability

Generative models and their physics genes



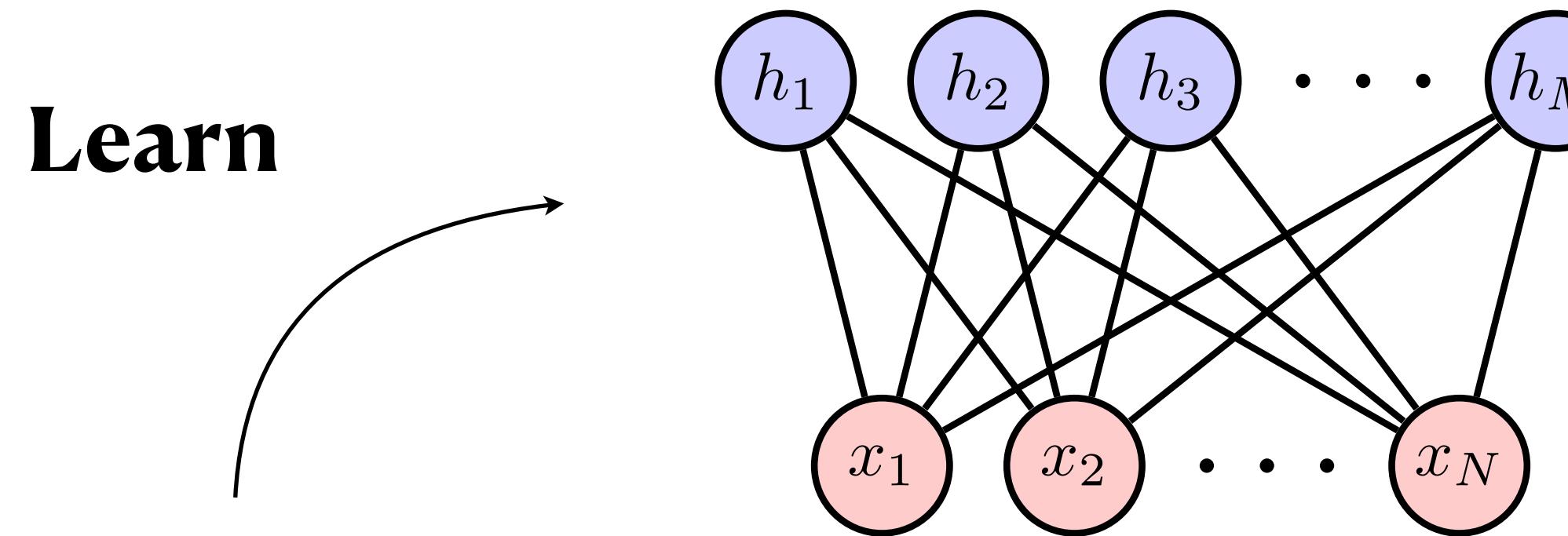
**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



Boltzmann machines

$$\mathcal{L} = -\mathbb{E}_{x \sim \text{data}} [\ln p(x)] \quad p(x) = e^{-E(x)}/Z$$



6	2	7	4	2	1	9
1	2	5	3	0	7	5
8	1	8	4	2	6	6
0	7	9	8	6	3	2
7	5	0	5	7	9	5
1	8	7	0	6	5	0
7	5	4	8	4	4	7

$$\nabla_{\theta} \mathcal{L} = \langle \nabla_{\theta} E \rangle_{\text{data}} - \langle \nabla_{\theta} E \rangle_{\text{model}}$$

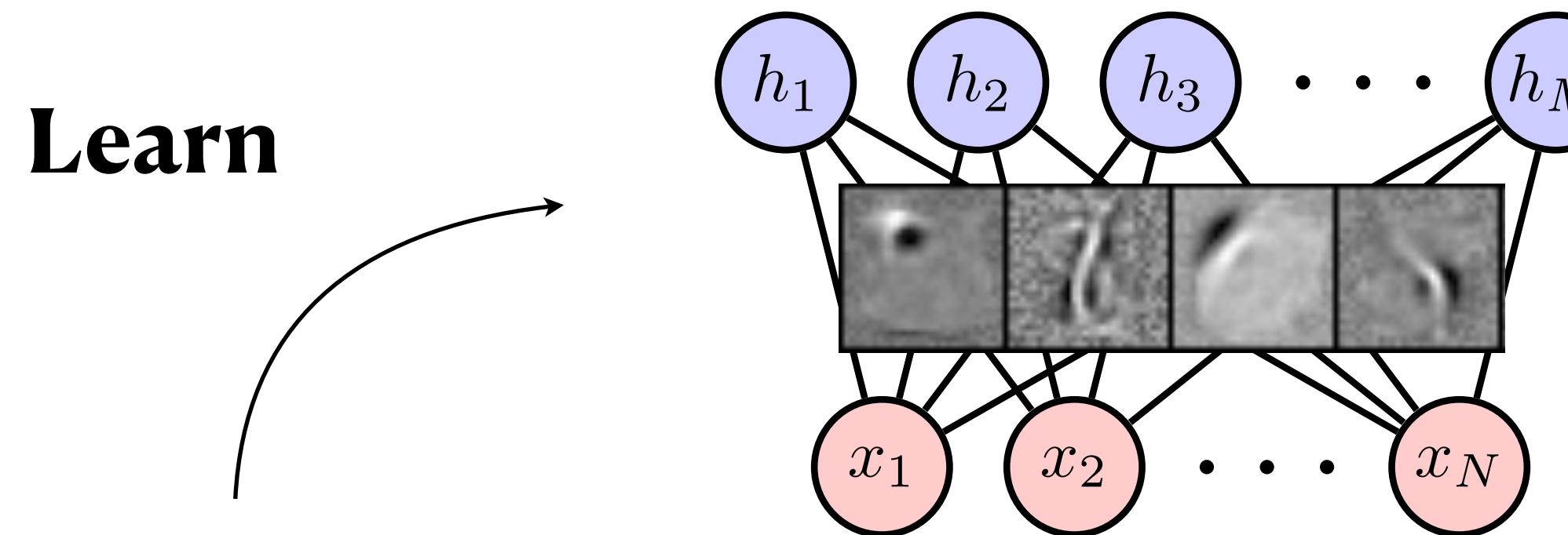
GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS

2210.10318

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

Boltzmann machines

$$\mathcal{L} = -\mathbb{E}_{x \sim \text{data}} [\ln p(x)] \quad p(x) = e^{-E(x)}/Z$$



6274219
1253045
8184266
0798630
7505795
1870650
7548447

$$\nabla_{\theta}\mathcal{L} = \langle \nabla_{\theta}E \rangle_{\text{data}} - \langle \nabla_{\theta}E \rangle_{\text{mode}}$$

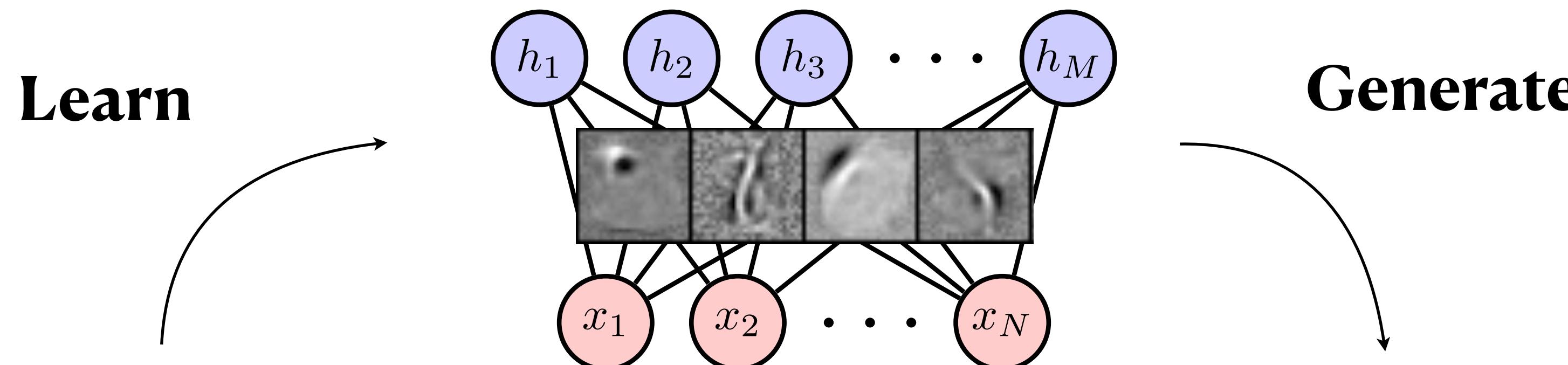
GAUSSIAN-BERNOULLI RBMS WITHOUT TEARS

2210.10318

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

Boltzmann machines

$$\mathcal{L} = -\mathbb{E}_{x \sim \text{data}} [\ln p(x)] \quad p(x) = e^{-E(x)}/Z$$



6 2 7 4 2 1 9
1 2 5 3 0 4 5
8 1 8 4 2 6 6
0 7 9 8 6 3 2
7 5 0 5 7 9 5
1 8 7 0 6 5 0
7 5 4 8 4 4 7

$$\nabla_{\theta}\mathcal{L} = \langle \nabla_{\theta}E \rangle_{\text{data}} - \langle \nabla_{\theta}E \rangle_{\text{model}}$$

1 8 3 1 6 7 1
6 6 3 3 3 6 8
4 5 8 4 4 1 9
3 7 7 9 8 7 6
1 5 3 5 0 2 2
4 2 5 1 2 4 2
3 0 5 0 7 0 9

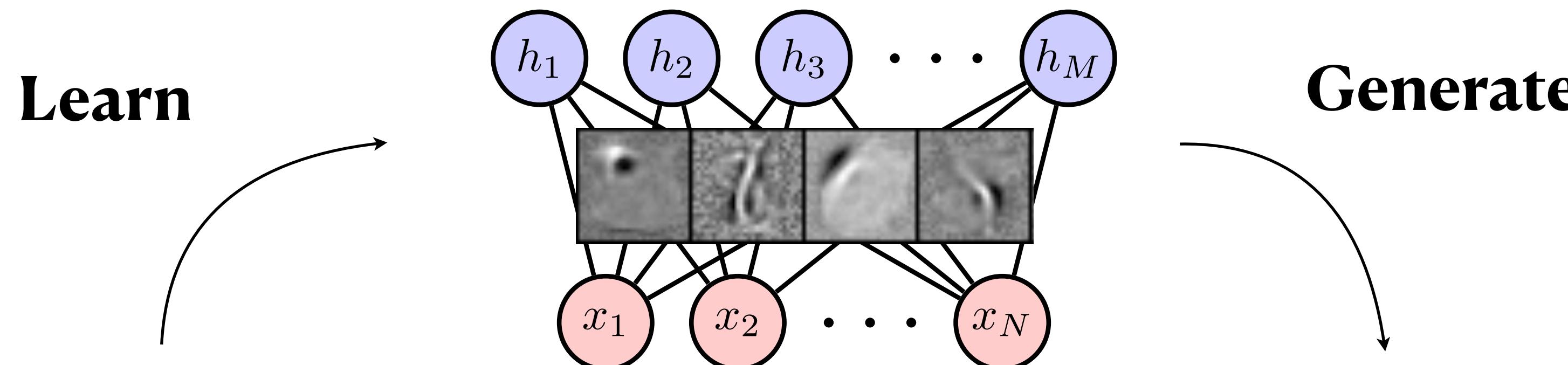
GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS

2210.10318

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

Boltzmann machines

$$\mathcal{L} = - \mathbb{E}_{x \sim \text{data}} [\ln p(x)] \quad p(x) = e^{-E(x)}/Z$$



6 2 7 4 2 1 9
1 2 5 3 0 7 5
8 1 8 4 2 6 6
0 7 9 8 6 3 2
7 5 0 5 7 9 5
1 8 7 0 6 5 0
7 5 4 8 4 4 7

$$\nabla_{\theta} \mathcal{L} = \langle \nabla_{\theta} E \rangle_{\text{data}} - \langle \nabla_{\theta} E \rangle_{\text{model}}$$

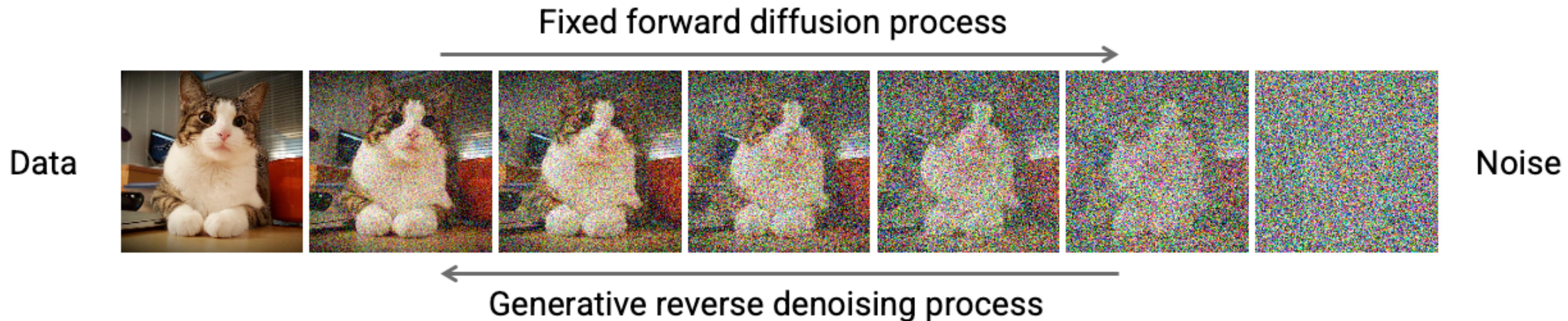
1 8 3 1 6 7 1
6 6 3 3 3 6 8
4 5 8 4 4 1 9
3 7 7 9 8 7 6
1 5 3 5 0 2 2
4 2 5 1 2 4 2
3 0 5 0 7 0 9

GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS

2210.10318

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

Diffusion models



I will follow the score-matching route <https://yang-song.net/blog/2021/score/>

Score matching

Minimizing Fisher divergence avoids the intractable partition function problem

$$\mathbb{F}(\pi \parallel p) = \int dx \pi(x) \left| \nabla_x \ln \pi(x) - \nabla_x \ln p(x) \right|^2$$

↑ ↑
target model

However, it brings up another problem

How to learn the model without knowing $\nabla \ln \pi$?

Implicit score matching

Integrate by parts Hyvarinen JMLR '05

$$\mathbb{F}(\pi \parallel p) = \int d\mathbf{x} \, \pi(\mathbf{x}) \left(|\nabla \ln p(\mathbf{x})|^2 + 2 \nabla^2 \ln p(\mathbf{x}) \right) + \text{const.}$$

The laplacian term can be difficult to compute

Cheaper stochastic estimate:
Song et al, 1905.07088

Curiously, the same expression for the kinetic energy of a wavefunction

$$\frac{\nabla^2 \psi}{\psi} = |\nabla \ln \psi|^2 + \nabla^2 \ln \psi$$

Forward laplacian:
Li et al, 2307.08214

Denoising score matching

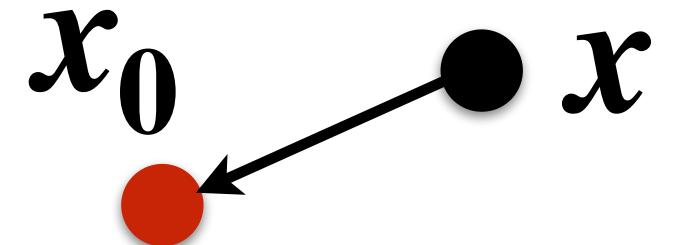
Perturb data with small noise Vincent 2011

$$q(x) = \int q(x | x_0) \pi(x_0) dx_0 \quad q(x | x_0) = \mathcal{N}(x; x_0, \sigma^2)$$

Fisher divergence between perturbed data and model is computable

$$\begin{aligned} \mathbb{F}(q \parallel p) &= \mathbb{E}_{x \sim q(x)} |\nabla \ln q(x) - \nabla \ln p(x)|^2 \\ &= \mathbb{E}_{x_0 \sim \pi(x_0)} \mathbb{E}_{x \sim q(x|x_0)} |\nabla \ln q(x | x_0) - \nabla \ln p(x)|^2 + \text{const.} \end{aligned}$$

$\hookrightarrow \frac{x_0 - x}{\sigma^2}$ the restoring force



Claim:

$$\mathbb{E}_{x \sim q(x)} |\nabla \ln q(x) - s_\theta|^2 = \mathbb{E}_{x_0 \sim \pi(x_0)} \mathbb{E}_{x \sim q(x|x_0)} |\nabla \ln q(x|x_0) - s_\theta|^2 + \text{const.}$$



$$\text{score } s_\theta = \nabla \ln p_\theta(x)$$

Independent
of θ

Proof:

$$\mathbb{E}_{x_0 \sim \pi(x_0)} \mathbb{E}_{x \sim q(x|x_0)} |s|^2 = \int dx_0 \int dx \pi(x_0) q(x|x_0) |s|^2 = \int dx q(x) |s|^2 = \mathbb{E}_{x \sim q(x)} |s|^2$$

$$\mathbb{E}_{x_0 \sim \pi(x_0)} \mathbb{E}_{x \sim q(x|x_0)} [s \cdot \nabla \ln q(x|x_0)] = \int dx_0 \int dx \pi(x_0) q(x|x_0) \frac{s \cdot \nabla q(x|x_0)}{q(x|x_0)}$$

$$= \int dx_0 \int dx \pi(x_0) s \cdot \nabla q(x|x_0)$$

$$= \int dx s \cdot \nabla q(x) = \mathbb{E}_{x \sim q(x)} [s \cdot \nabla \ln q(x)]$$

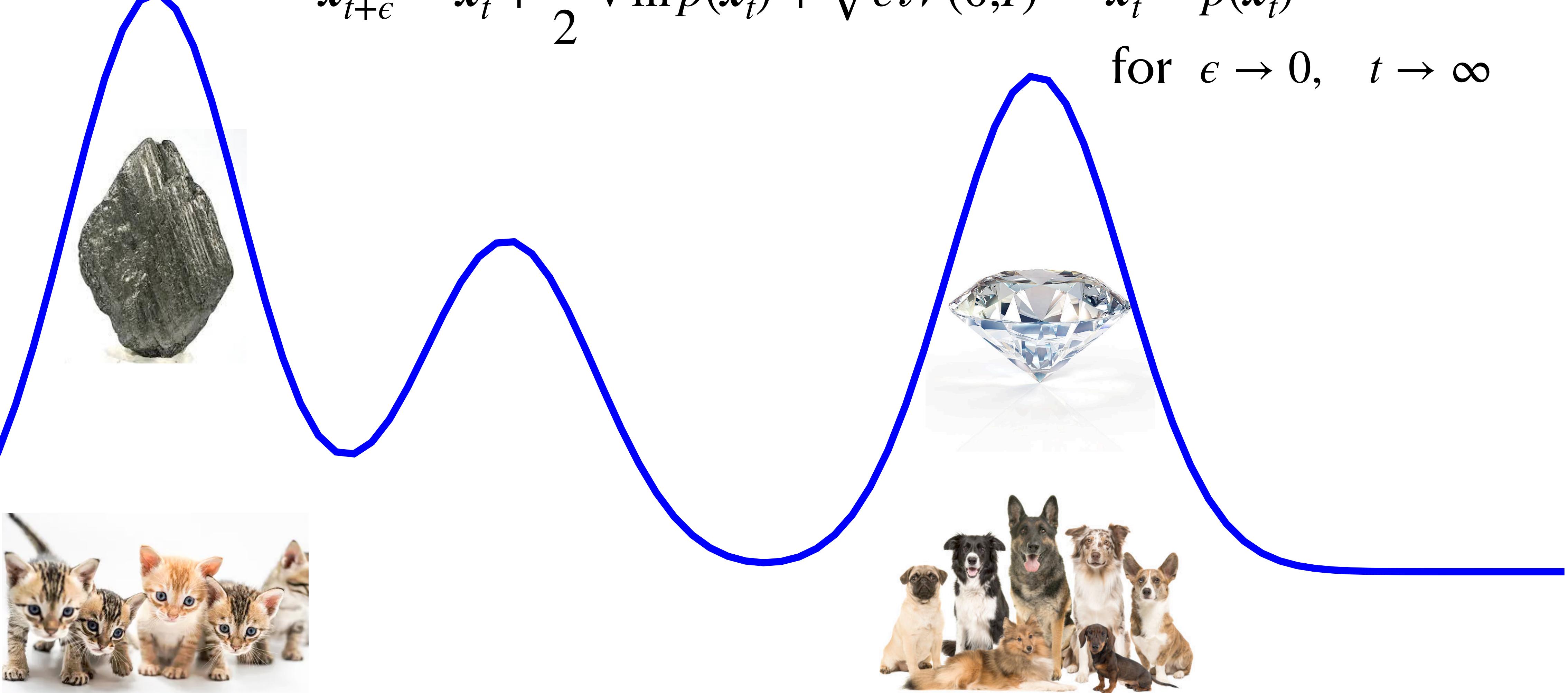
Why score matching did not take off?

Hard to sample between modes with Langevin dynamics

$$x_{t+\epsilon} = x_t + \frac{\epsilon}{2} \nabla \ln p(x_t) + \sqrt{\epsilon} \mathcal{N}(0, I)$$

$$x_t \sim p(x_t)$$

for $\epsilon \rightarrow 0, t \rightarrow \infty$



From denoising score matching to diffusion model

Song et al, Generative modeling by estimating gradients of the data distribution, 1907.05600

Built upon this intuition, we propose to improve score-based generative modeling by 1) *perturbing the data using various levels of noise*; and 2) *simultaneously estimating scores corresponding to all noise levels by training a single conditional score network*. After training, when using Langevin dynamics to generate samples, we initially use scores corresponding to large noise, and gradually anneal down the noise level. This helps smoothly transfer the benefits of large noise levels to low noise levels where the perturbed data are almost indistinguishable from the original ones. In what follows, we will elaborate more on the details of our method, including the architecture of our score networks, the training objective, and the annealing schedule for Langevin dynamics.

Sohl-Dickstein et al, Deep unsupervised learning using nonequilibrium thermodynamics, 1503.03585

The essential idea, inspired by non-equilibrium statistical physics, is to systematically and slowly destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This ap-

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.

From denoising score matching to diffusion model

The objective of denoising diffusion probabilistic model

Song et al, 1907.05600
Ho et al, 2006.11239

<https://cvpr2022-tutorial-diffusion-models.github.io>

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)} \| \mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0) \|_2^2$$

diffusion time t data sample \mathbf{x}_0 diffused data sample \mathbf{x}_t neural network score of diffused data sample

Sample with **annealed Langevin dynamics** with decreasing steps ϵ_t

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\epsilon_t}{2} \mathbf{s}(\mathbf{x}_t, t) + \sqrt{\epsilon_t} \mathcal{N}(0, I)$$

A tale of three equations

Langevin equation (SDE)

$$\mathbf{x}_{t+dt} = \mathbf{x}_t + f dt + \sqrt{2dt} \mathcal{N}(0, I)$$

Fokker-Planck equation (PDE)

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla \cdot [p(\mathbf{x}, t) \mathbf{f}] - \nabla^2 p(\mathbf{x}, t) = 0$$

“Particle method” (ODE)

$$\frac{d\mathbf{x}}{dt} = \mathbf{f} - \nabla \ln p(\mathbf{x}, t) \equiv \mathbf{v}$$

(Another way to reverse the diffusion is
via the reverse-time SDE Anderson 1982)

Maoutsou et al, 2006.00702
Song et al, 2011.13456

A tale of three equations

Langevin equation (SDE)

$$\mathbf{x}_{t+dt} = \mathbf{x}_t + \mathbf{f} dt + \sqrt{2dt} \mathcal{N}(0, I)$$

Fokker-Planck equation (PDE)

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla \cdot \left[p(\mathbf{x}, t) (\mathbf{f} - \nabla \ln p(\mathbf{x}, t)) \right] = 0$$

“Particle method” (ODE)

$$\frac{d\mathbf{x}}{dt} = \mathbf{f} - \nabla \ln p(\mathbf{x}, t) \equiv \mathbf{v}$$

(Another way to reverse the diffusion is
via the reverse-time SDE Anderson 1982)

Maoutsou et al, 2006.00702
Song et al, 2011.13456

$$\mathcal{P}(\vec{x}, t) = \int d^3\vec{x}' \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} \exp \left[-\frac{(\vec{x} - \vec{x}' - \epsilon \vec{v}(\vec{x}'))^2}{4D\epsilon} \right] \mathcal{P}(\vec{x}', t - \epsilon), \quad (9.18)$$

and simplified by the change of variables,

$$\begin{aligned} \vec{y} &= \vec{x}' + \epsilon \vec{v}(\vec{x}') - \vec{x} \implies \\ d^3\vec{y} &= d^3\vec{x}' (1 + \epsilon \nabla \cdot \vec{v}(\vec{x}')) = d^3\vec{x}' (1 + \epsilon \nabla \cdot \vec{v}(\vec{x}) + \mathcal{O}(\epsilon^2)). \end{aligned} \quad (9.19)$$

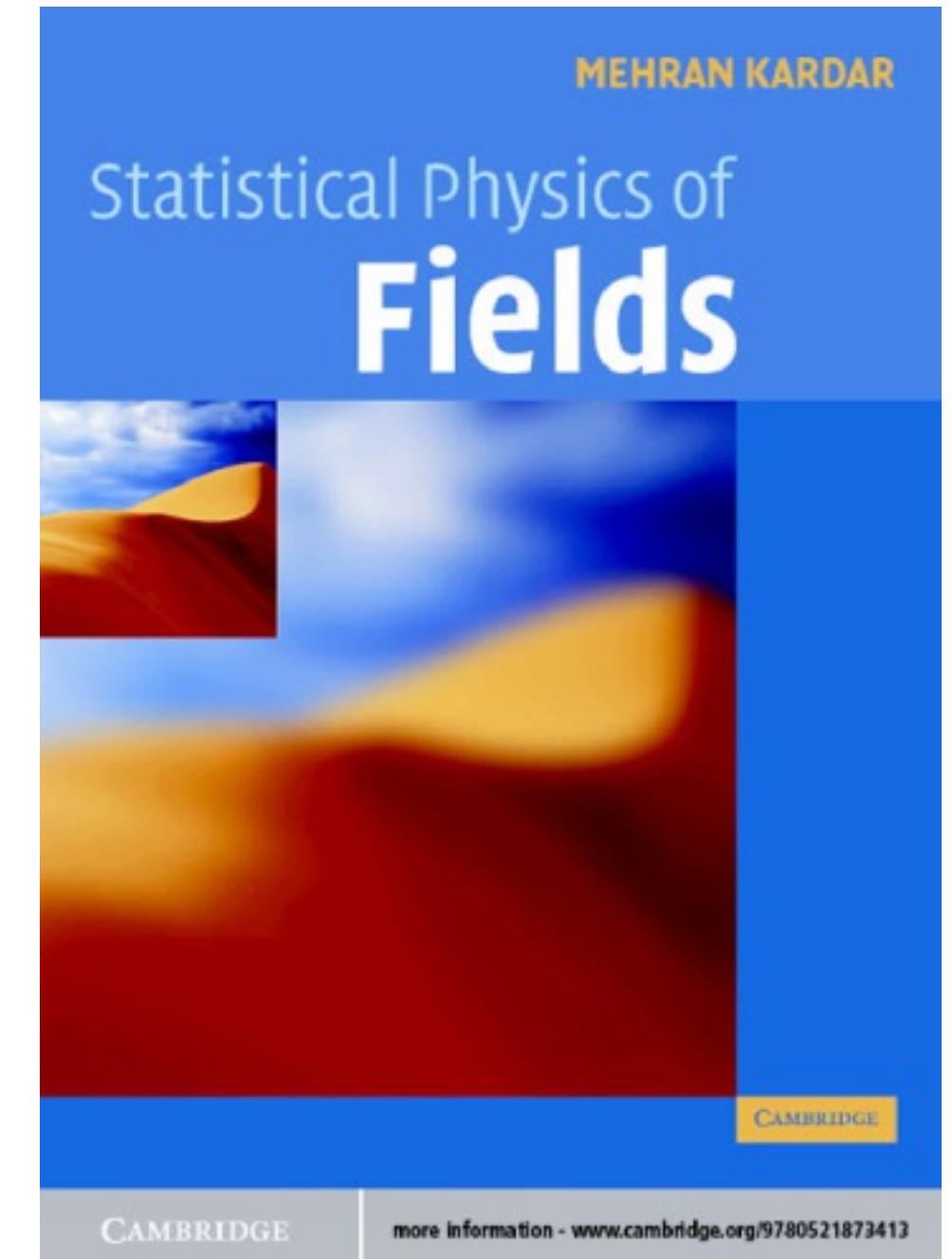
Keeping only terms at order of ϵ , we obtain

$$\begin{aligned} \mathcal{P}(\vec{x}, t) &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \int d^3\vec{y} \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} e^{-\frac{y^2}{4D\epsilon}} \mathcal{P}(\vec{x} + \vec{y} - \epsilon \vec{v}(\vec{x}), t - \epsilon) \\ &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \int d^3\vec{y} \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} e^{-\frac{y^2}{4D\epsilon}} \\ &\times \left[\mathcal{P}(\vec{x}, t) + (\vec{y} - \epsilon \vec{v}(\vec{x})) \cdot \nabla \mathcal{P} + \frac{y_i y_j - 2\epsilon y_i v_j + \epsilon^2 v_i v_j}{2} \nabla_i \nabla_j \mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^2) \right] \\ &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \left[\mathcal{P} - \epsilon \vec{v} \cdot \nabla + \epsilon D \nabla^2 \mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^2) \right]. \end{aligned} \quad (9.20)$$

Equating terms at order of ϵ leads to the *Fokker–Planck equation*,

$$\frac{\partial \mathcal{P}}{\partial t} + \nabla \cdot \vec{J} = 0, \quad \text{with} \quad \vec{J} = \vec{v} \mathcal{P} - D \nabla \mathcal{P}. \quad (9.21)$$

from Langevin
to Fokker-Planck



Lessons from diffusion models

Continuous normalizing flow has great potential: diffusion model is an “existence proof”

Going beyond maximum likelihood estimation training (even if we can)

[https://blog.alexalemi.com/
diffusion.html](https://blog.alexalemi.com/diffusion.html)

The conditional trick (originated from denoising score matching Vincent 2011)

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)} \| \mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0) \|_2^2$$

diffusion time t data sample \mathbf{x}_0 diffused data sample \mathbf{x}_t neural network score of diffused data sample

Lessons from diffusion models

Continuous normalizing flow has great potential: diffusion model is an “existence proof”

Going beyond maximum likelihood estimation training (even if we can)

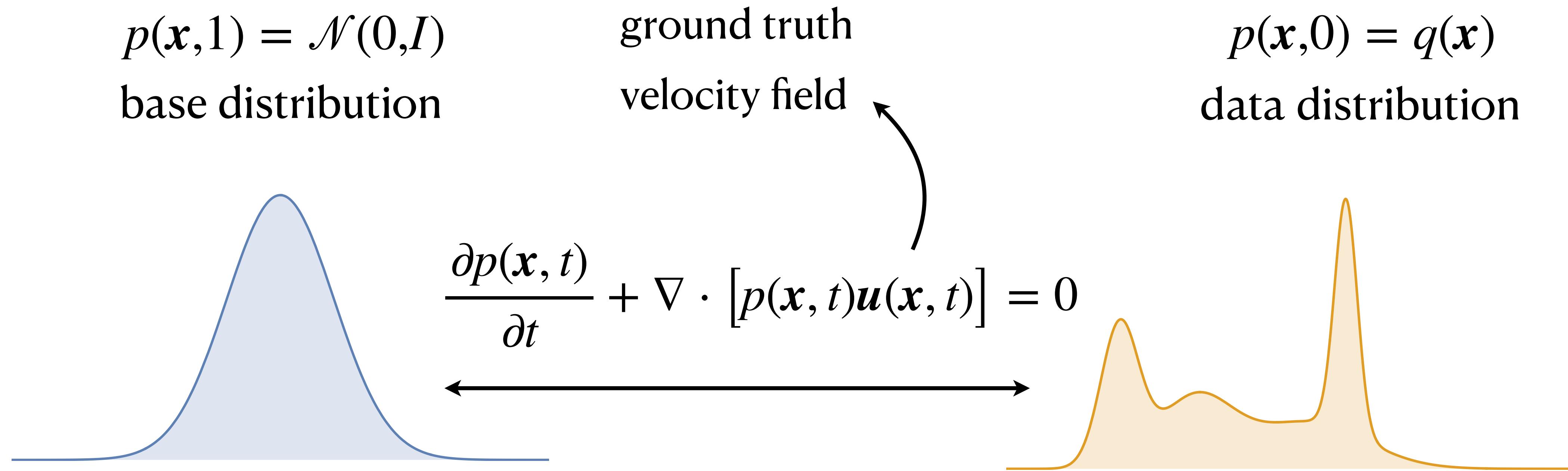
[https://blog.alexalemi.com/
diffusion.html](https://blog.alexalemi.com/diffusion.html)

The conditional trick (originated from denoising score matching Vincent 2011)

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t)} \| \mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \|_2^2$$

diffusion time t diffused data \mathbf{x}_t neural network score of diffused data (marginal)

Flow matching



$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}, t)} \left| \mathbf{v}_\theta(\mathbf{x}, t) - \mathbf{u}(\mathbf{x}, t) \right|^2$$

The “conditional” trick

Given a conditional continuity equation

$$\frac{\partial p(\mathbf{x} | \mathbf{x}_0, t)}{\partial t} + \nabla \cdot [p(\mathbf{x} | \mathbf{x}_0, t) \mathbf{u}(\mathbf{x} | \mathbf{x}_0, t)] = 0$$

Then, up to a constant, we have

$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x} | \mathbf{x}_0, t)} \left| \mathbf{v}_\theta(\mathbf{x}, t) - \mathbf{u}(\mathbf{x} | \mathbf{x}_0, t) \right|^2$$

We can learn the ground truth velocity by regressing on the conditional velocity

Claim: $\mathcal{L}_{\text{CFM}} = \mathcal{L}_{\text{FM}} + \text{const.}$

where $\mathcal{L}_{\text{FM}} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{x \sim p(x,t)} |v_\theta(x, t) - u(x, t)|^2$

$$\mathcal{L}_{\text{CFM}} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{x_0 \sim q(x_0)} \mathbb{E}_{x \sim p(x|x_0, t)} |v_\theta(x, t) - u(x|x_0, t)|^2$$

$$p(x, t) = \int p(x|x_0, t) q(x_0) dx_0 \quad p(x, t)u(x, t) = \int p(x|x_0, t)u(x|x_0, t) q(x_0) dx_0$$

Proof:

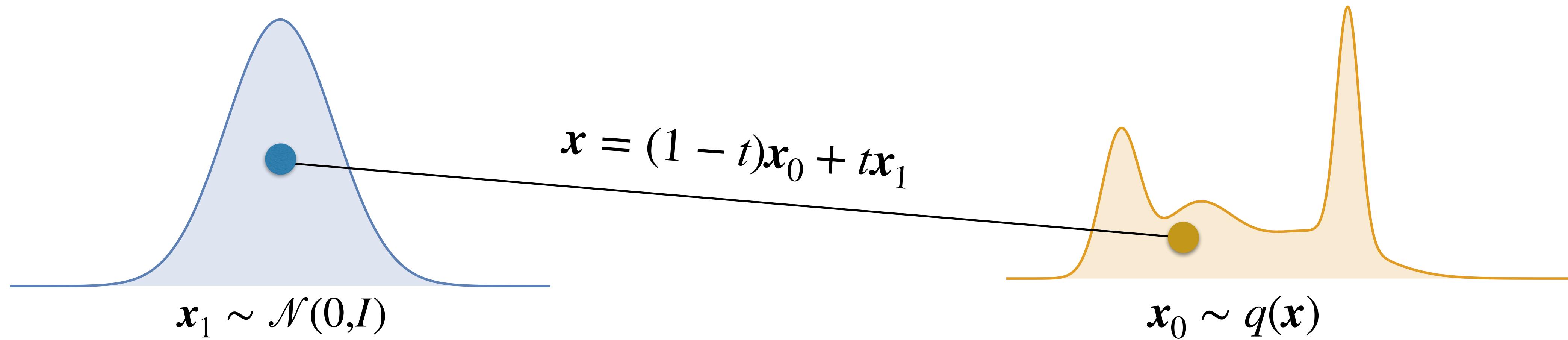
$$\mathbb{E}_{x_0 \sim q(x_0)} \mathbb{E}_{x \sim p(x|x_0, t)} |v_\theta|^2 = \int dx_0 \int dx q(x_0) p(x|x_0, t) |v_\theta|^2 = \int dx p(x, t) |v_\theta|^2 = \mathbb{E}_{x \sim p(x, t)} |v_\theta|^2$$

$$\mathbb{E}_{x_0 \sim q(x_0)} \mathbb{E}_{x \sim p(x|x_0, t)} [v_\theta \cdot u(x|x_0, t)] = \int dx_0 \int dx q(x_0) p(x|x_0, t) [v_\theta \cdot u(x|x_0, t)]$$

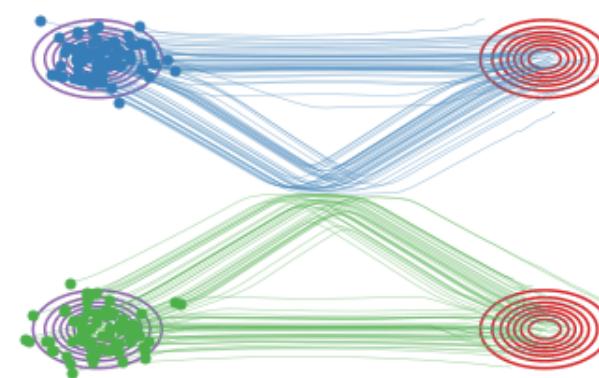
$$= \int dx p(x, t) v_\theta \cdot u(x, t) = \mathbb{E}_{x \sim p(x, t)} [v_\theta \cdot u(x, t)]$$

Examples of flow matching

$$p(x | x_0, t) = \mathcal{N}((1 - t)x_0, t^2) \quad u(x | x_0, t) = \frac{dx}{dt} = x_1 - x_0$$



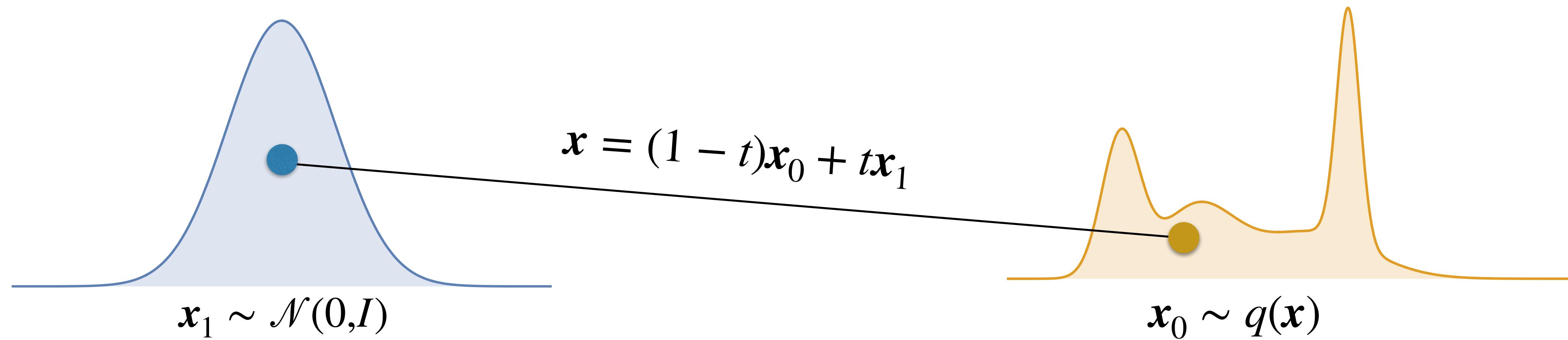
$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{x_0 \sim q(x_0)} \mathbb{E}_{x_1 \sim \mathcal{N}(0,I)} \left| v_\theta(x, t) - (x_1 - x_0) \right|^2$$



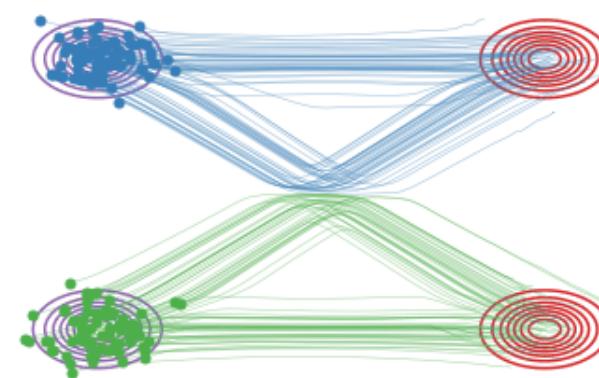
Causalizing linear interpolation with rectified flow 2209.03003
<https://www.cs.utexas.edu/~lqiang/rectflow/html/intro.html>

Examples of flow matching

$$p(x | x_0, t) = \mathcal{N}((1 - t)x_0, t^2) \quad u(x | x_0, t) = \frac{dx}{dt} = x_1 - x_0$$



$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{x_0 \sim q(x_0)} \mathbb{E}_{x_1 \sim \mathcal{N}(0,I)} \left| v_\theta(x, t) - (x_1 - x_0) \right|^2$$



Causalizing linear interpolation with rectified flow 2209.03003
<https://www.cs.utexas.edu/~lqiang/rectflow/html/intro.html>

Flow matching is all you need!

This framework contains various diffusion models as special cases

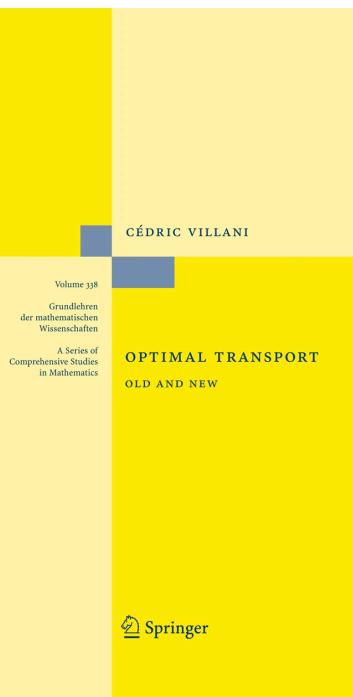
The base distribution does not have to be Gaussian

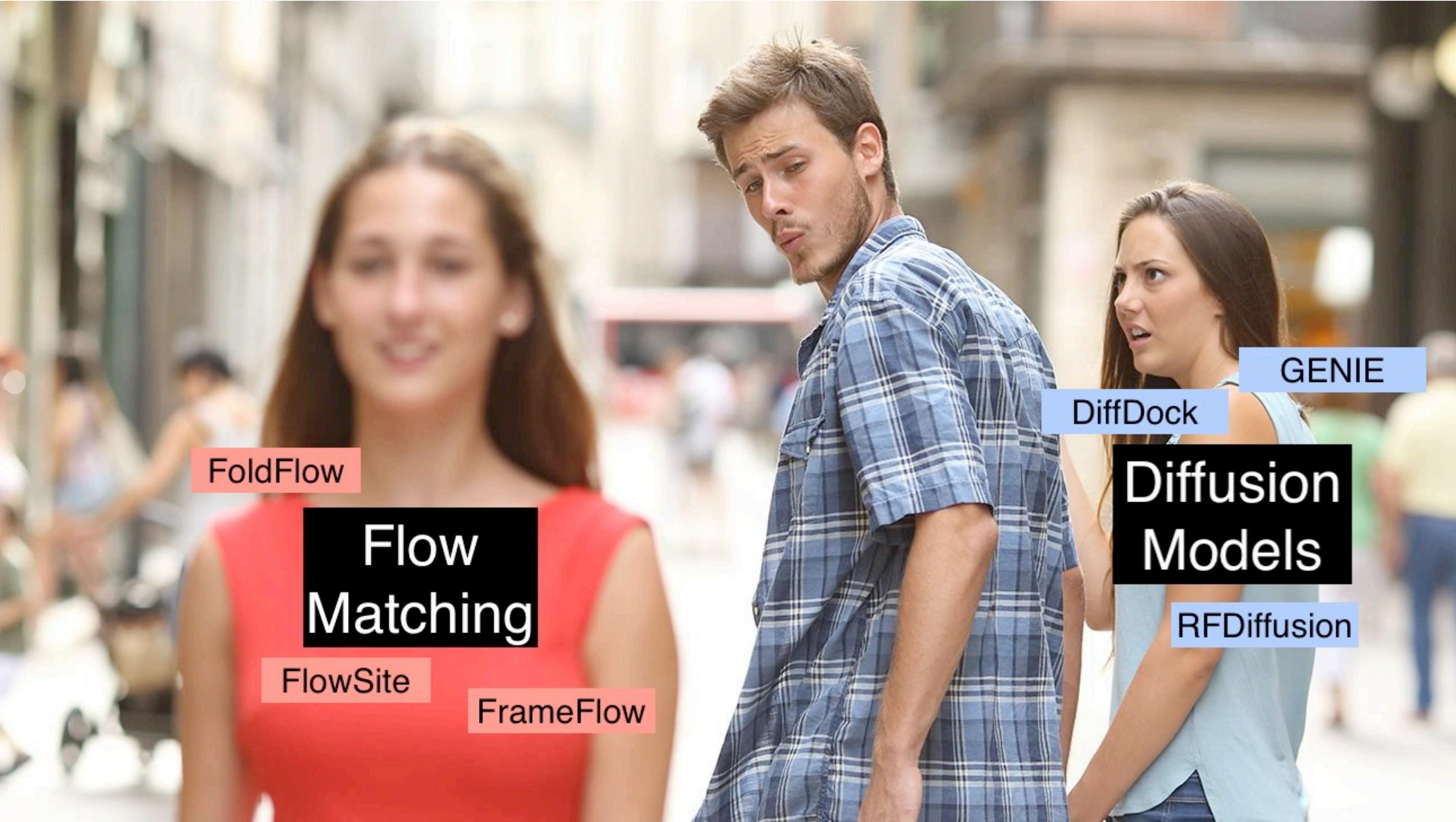
Fast generation with rectified transportation path (Liu et al 2209.03003)

400x speedup compared to continuous normalizing flow (Albergo et al, 2209.15571)

Surpasses diffusion model on Imagenet in likelihood and sample quality
(Lipman et al, 2210.02747)

Generalization to flow on Riemannian manifolds (Chen et al, 2302.03660)



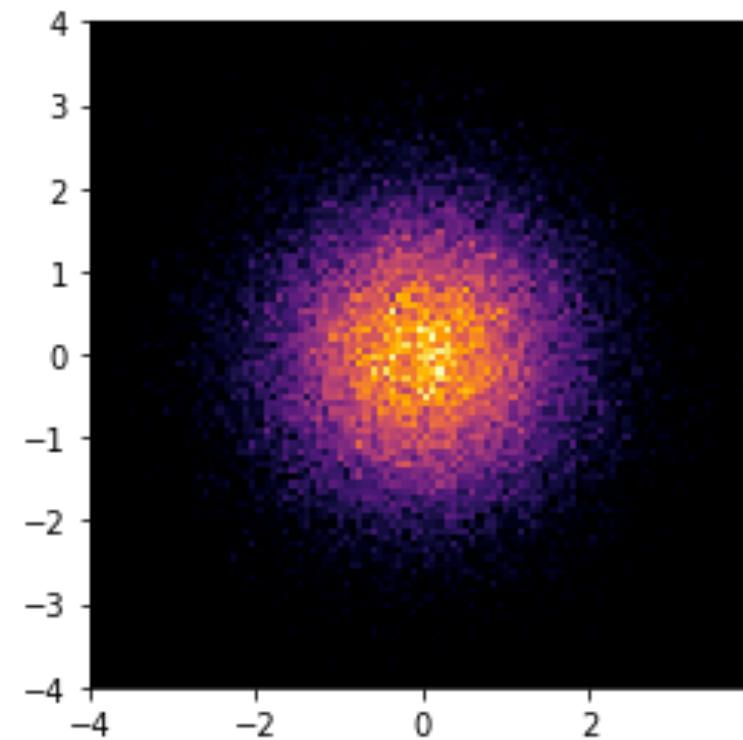


https://twitter.com/michael_galkin/status/1711845455817261409

Demo: free energy of classical Coulomb gas

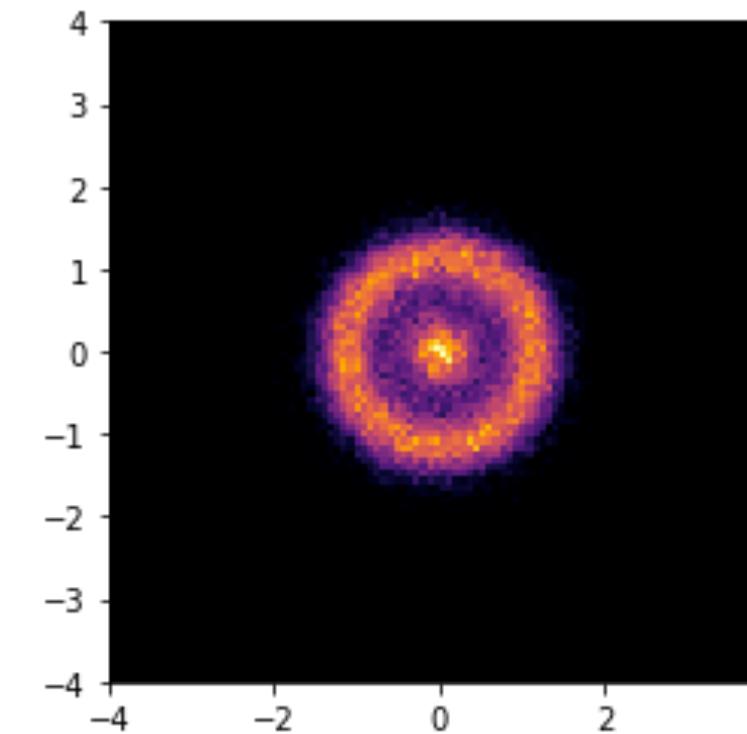
$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{x_0 \sim \mathcal{N}(0,I)} \mathbb{E}_{x_1 \sim \exp(-\beta E)/Z} \left| x_1 - x_0 - v(x, t) \right|^2$$

$$Z = \mathbb{E}_{x \sim q(x)} [e^{-\beta E(x) - \ln q(x)}] \quad \ln q(x) = \ln \mathcal{N}(0,I) - \int_0^1 \nabla \cdot v dt$$



Base density
Gaussian samples

←→
Interpolate samples to
estimate free energy
differences

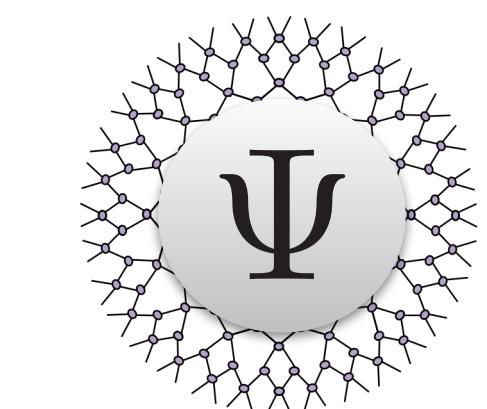


Target density
Monte Carlo samples



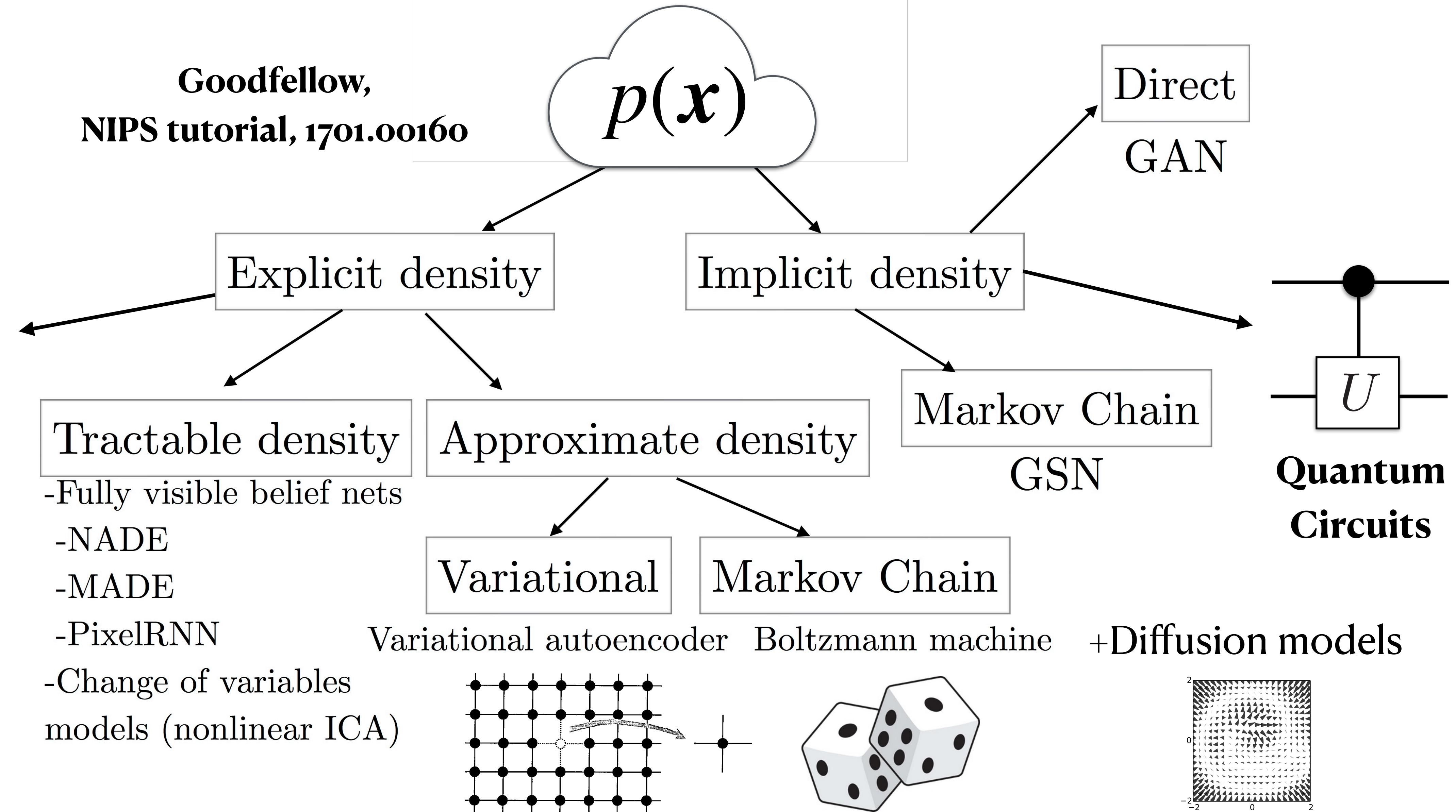
<https://colab.research.google.com/drive/1t-Vk37Axxpo4oB7uXFUNlk-zeCC2lcX3?usp=sharing>

Generative models and their physics genes



**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



Variational autoencoders

Kingma, Welling, 1312.6114

Close connection to the variational calculus we have learned

$$p(\mathbf{x}) = \frac{e^{-\beta E(\mathbf{x})}}{Z}$$

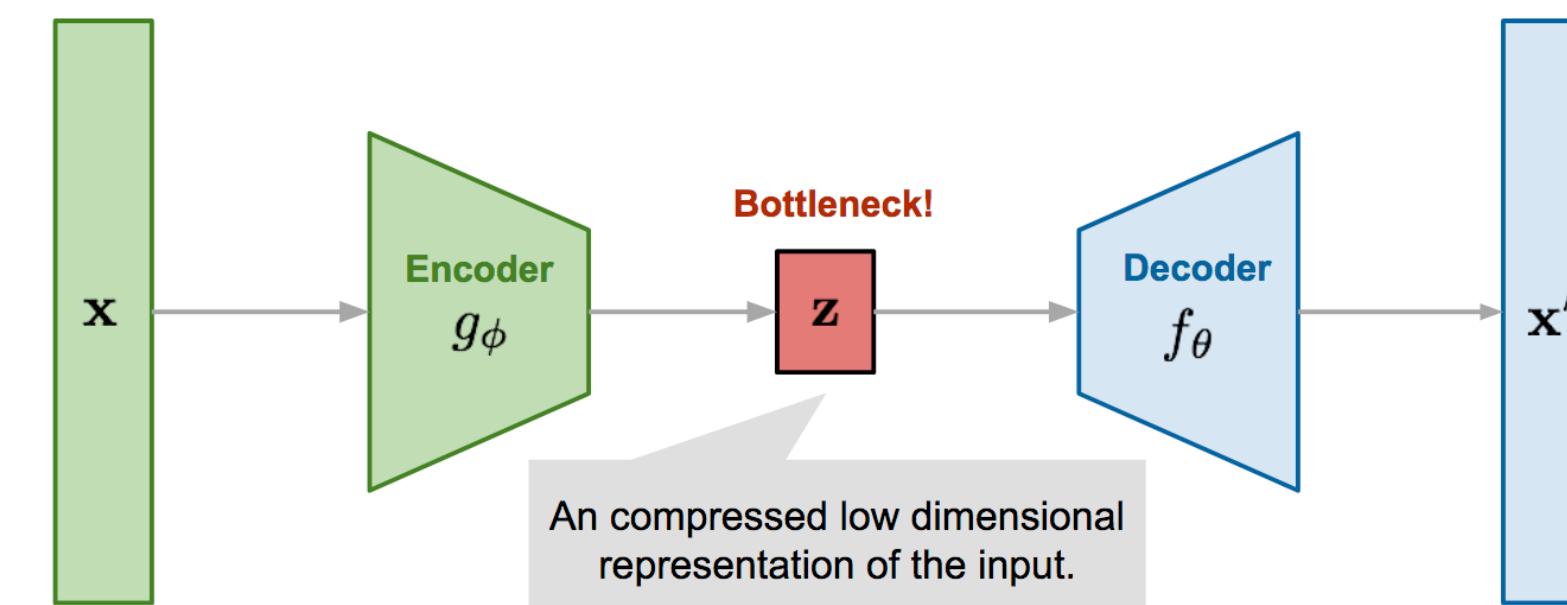
Variational free energy

$$\int d\mathbf{x} q(\mathbf{x}) [\ln q(\mathbf{x}) + \beta E(\mathbf{x})] \geq -\ln Z$$

$$p(z | \mathbf{x}) = \frac{p(\mathbf{x}, z)}{p(\mathbf{x})}$$

Variational Bayes/Variational inference

$$\int dz q(z | \mathbf{x}) [\ln q(z | \mathbf{x}) - \ln p(\mathbf{x}, z)] \geq -\ln p(\mathbf{x})$$



For each data we introduce

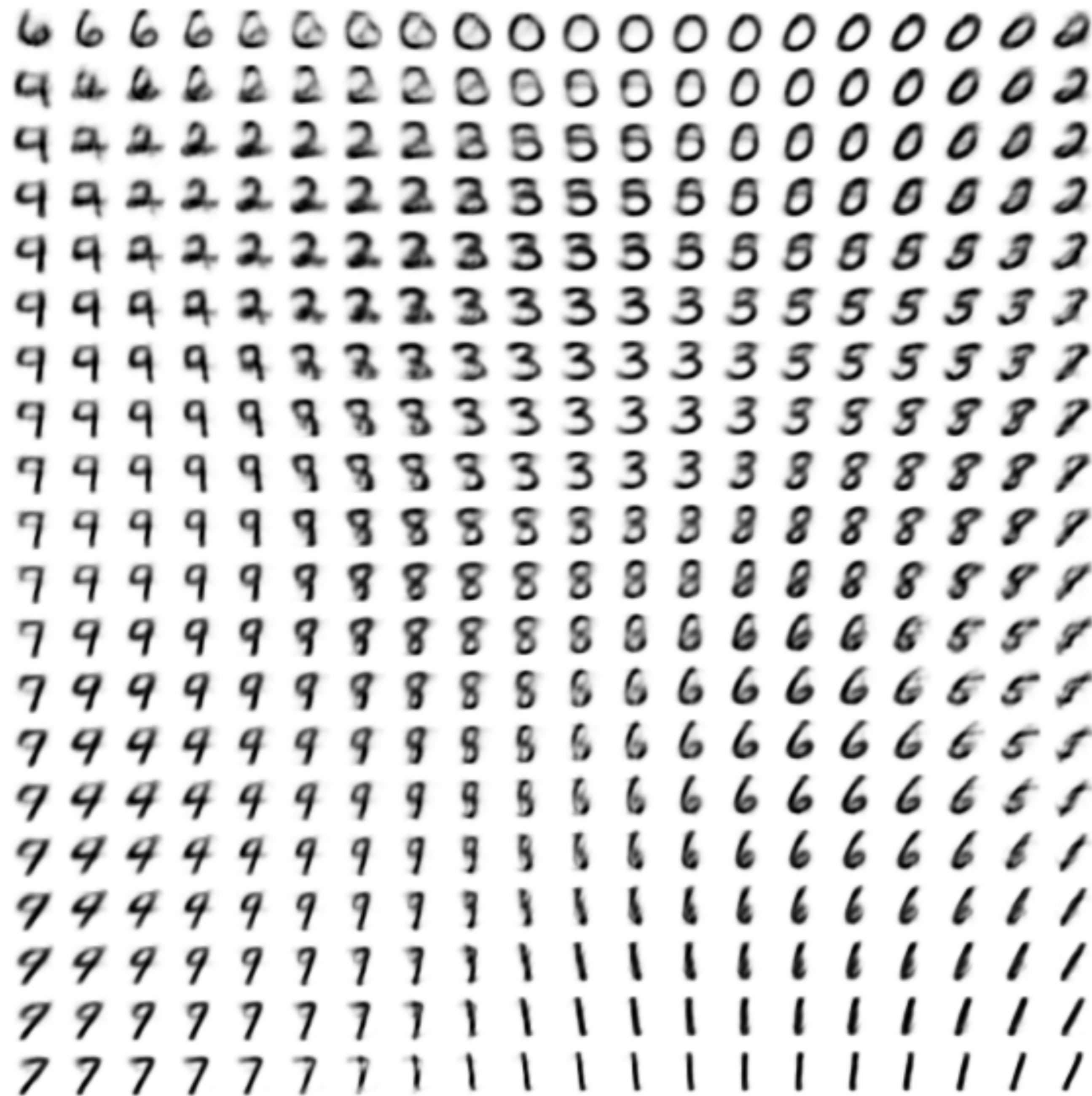
$$\mathcal{L}(\mathbf{x}) = \langle -\ln p(\mathbf{x}, \mathbf{z}) + \ln q(\mathbf{z}|\mathbf{x}) \rangle_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})}, \quad (53)$$

which is a variational upper bound of $-\ln p(\mathbf{x})$ since $\mathcal{L}(\mathbf{x}) + \ln p(\mathbf{x}) = \text{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \geq 0$. We see that $q(\mathbf{z}|\mathbf{x})$ provides a variational approximation of the posterior $p(\mathbf{z}|\mathbf{x})$. By minimizing \mathcal{L} one effectively pushes the two distributions together. And the variational free energy becomes exact only when $q(\mathbf{z}|\mathbf{x})$ matches to $p(\mathbf{z}|\mathbf{x})$. In fact, $-\mathcal{L}$ is called evidence lower bound (ELBO) in variational inference.

We can obtain an alternative form of the variational free energy

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = -\langle \ln p_{\theta}(\mathbf{x}|\mathbf{z}) \rangle_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} + \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})). \quad (54)$$

The first term of Eq. (54) is the reconstruction negative log-likelihood, while the second term is the KL divergence between the approximate posterior distribution and the latent prior. We also be explicit about the network parameters θ, ϕ of the encoder and decoder.

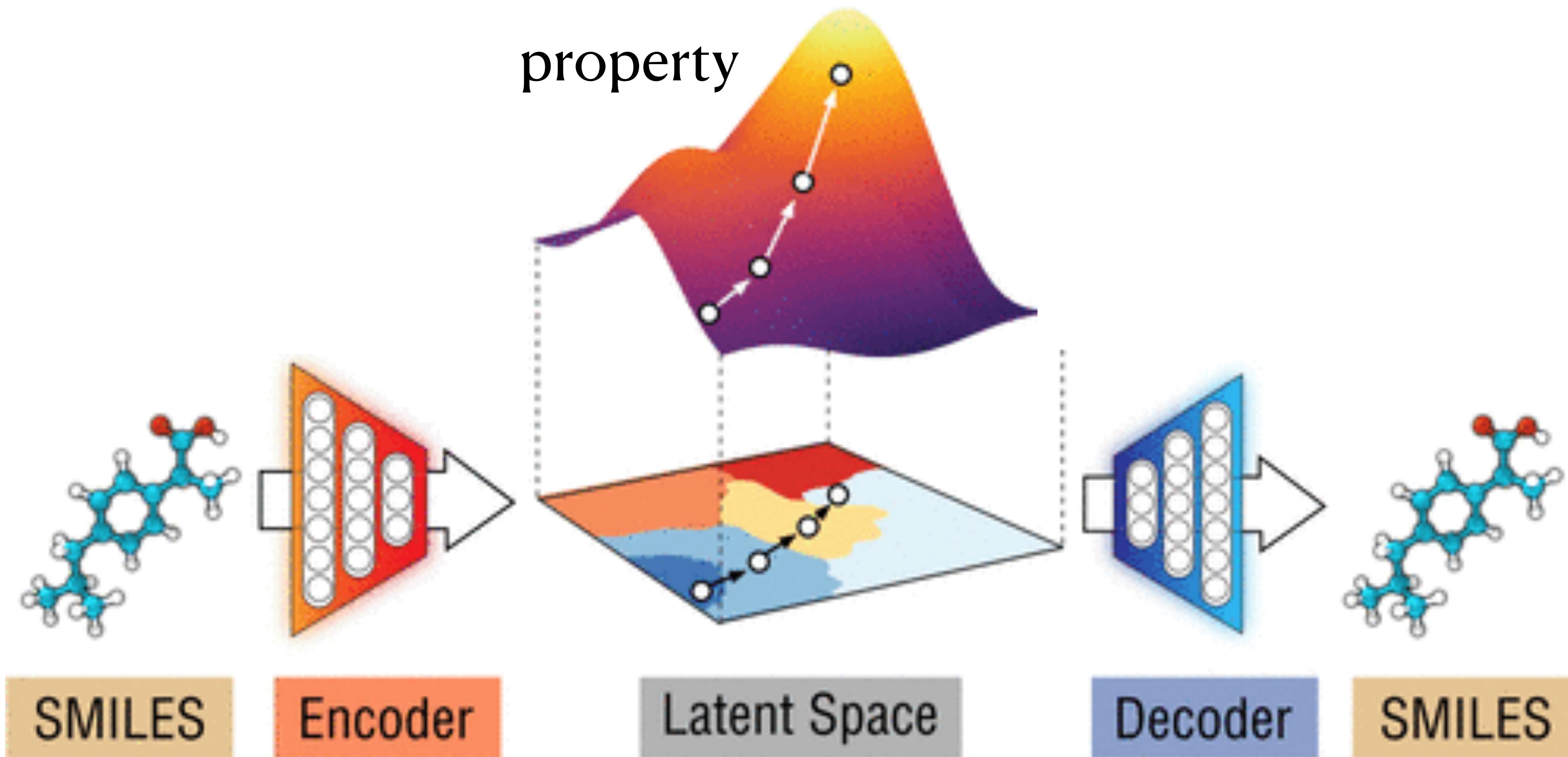


Learned MNIST latent space

Kingma, Welling, 1312.6114

Chemical design using continuous latent variables

Gomez-Bombarelli et al, 1610.02415



GAN

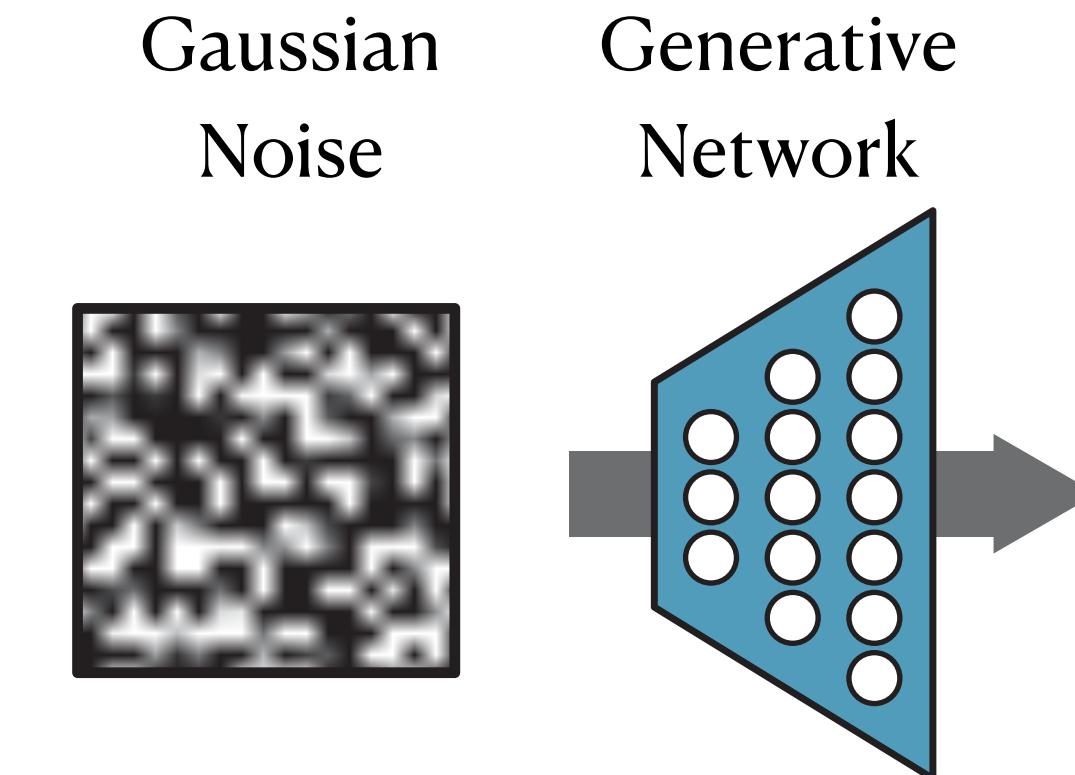
Likelihood free simulator

Prone to mode collapse “de”generate

More tricky to train than others

Performance have been surpassed by diffusion models

I found GAN to be less useful for quantitative scientific applications



<https://www.christies.com/Features/A-collaboration-between-two-artists-one-human-one-a-machine-9332-1.aspx>

GAN

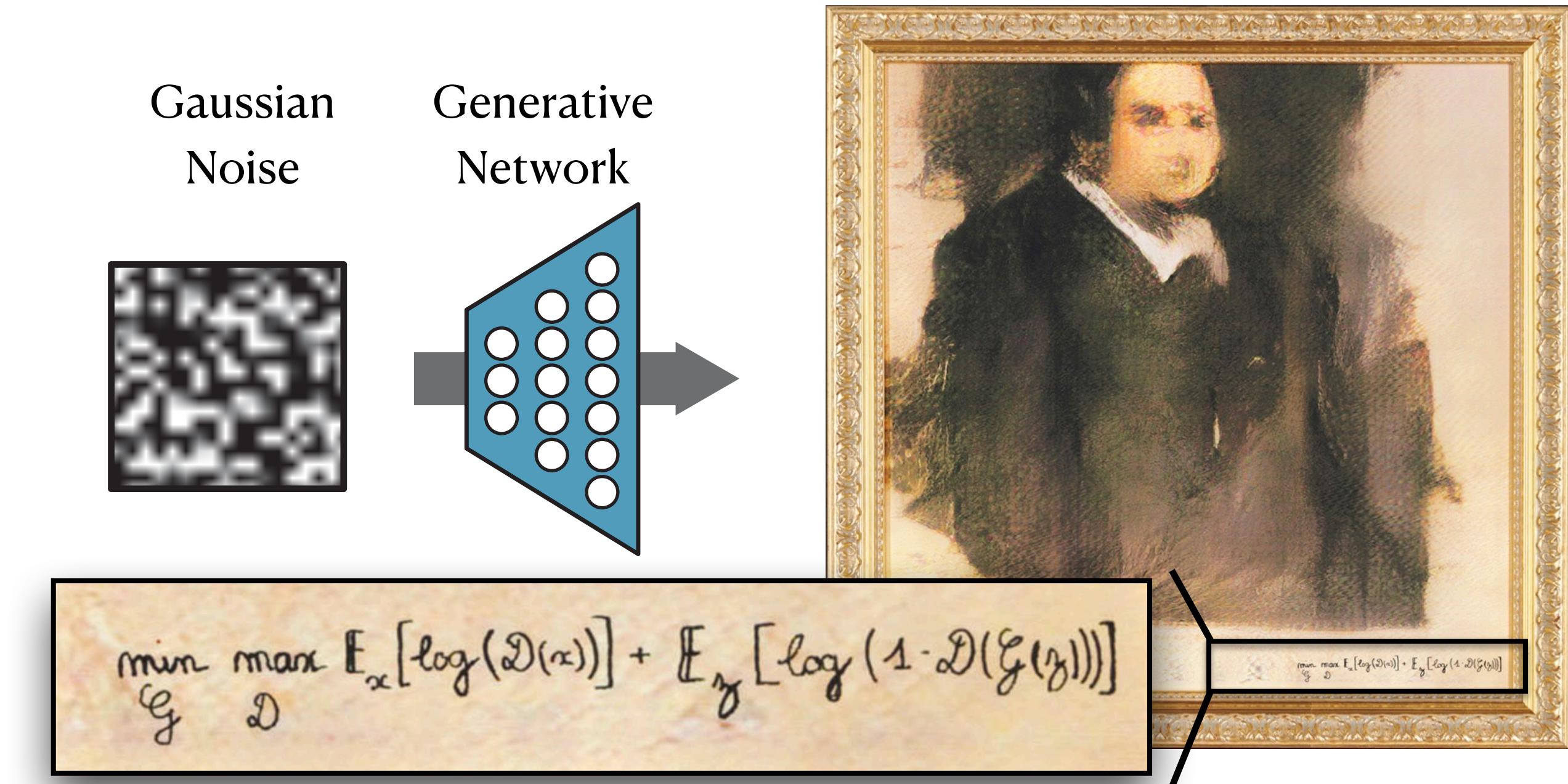
Likelihood free simulator

Prone to mode collapse “de”generate

More tricky to train than others

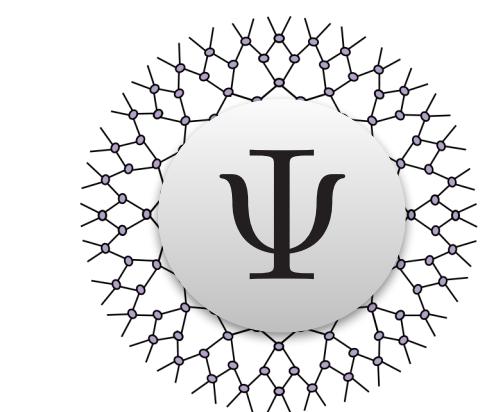
Performance have been surpassed by diffusion models

I found GAN to be less useful for quantitative scientific applications



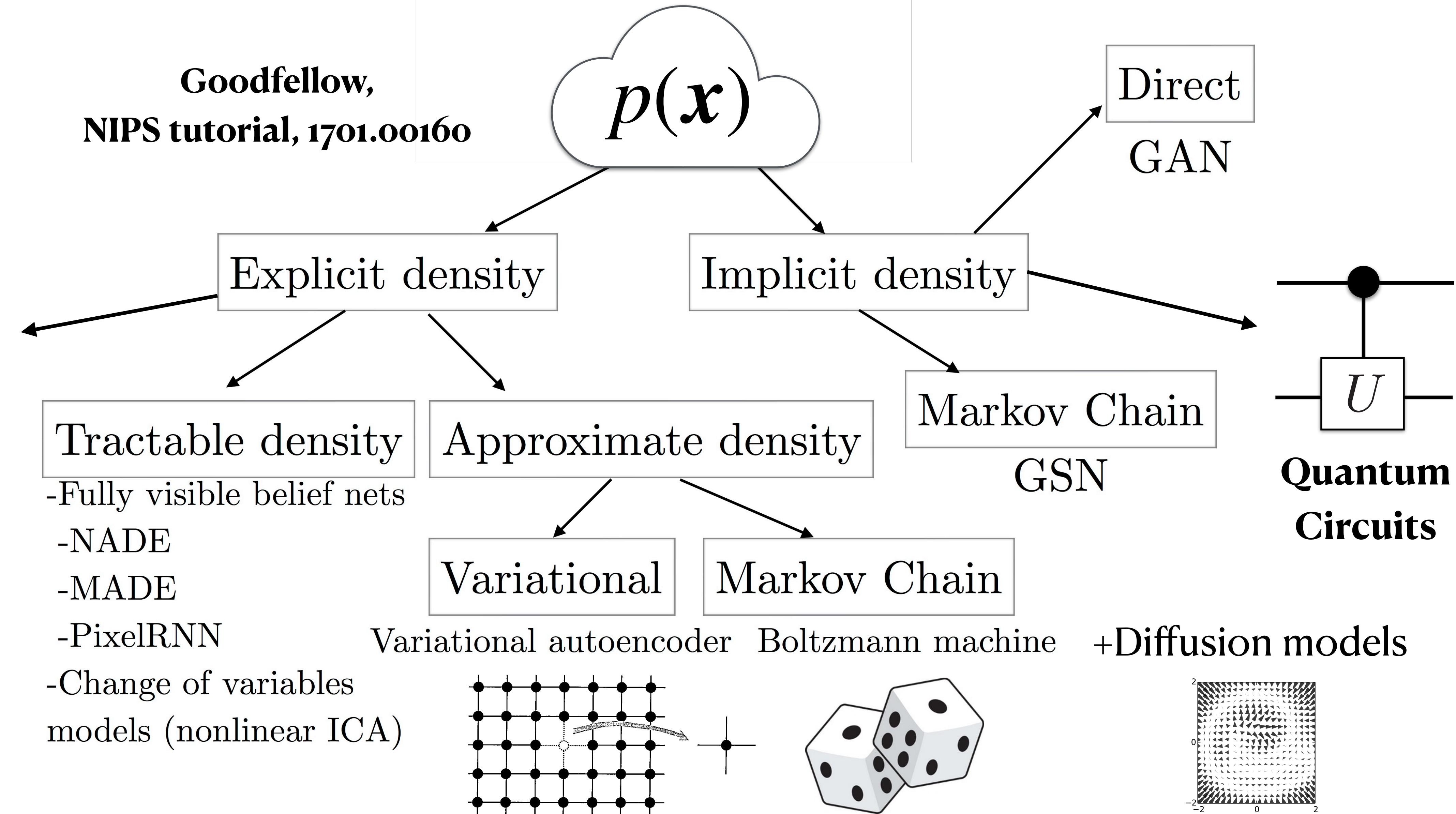
<https://www.christies.com/Features/A-collaboration-between-two-artists-one-human-one-a-machine-9332-1.aspx>

Generative models and their physics genes



**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



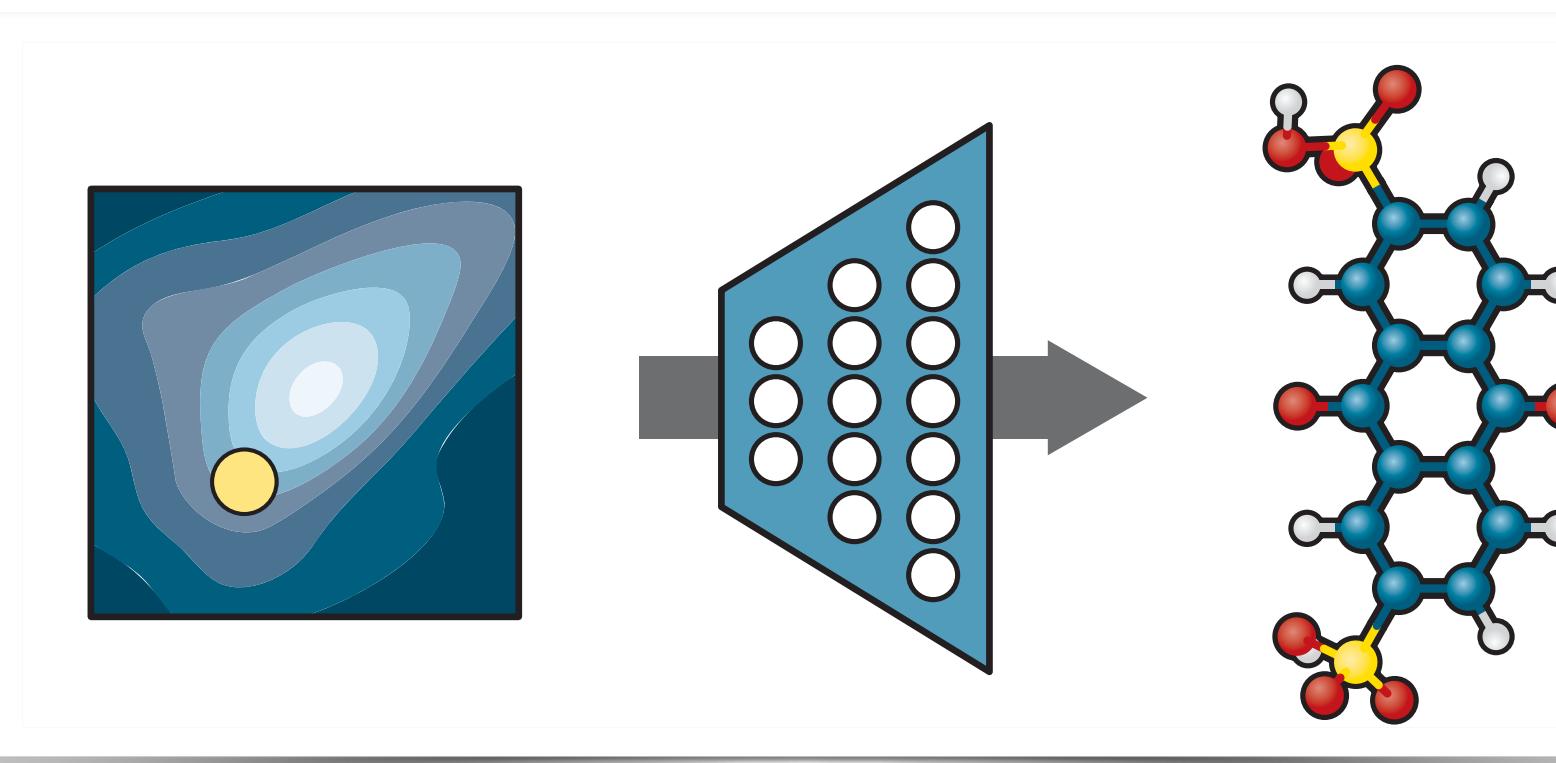
Generative AI for Science

1

How to Build a GPT-3 for
Science

Scientific language model

2



Matter inverse design

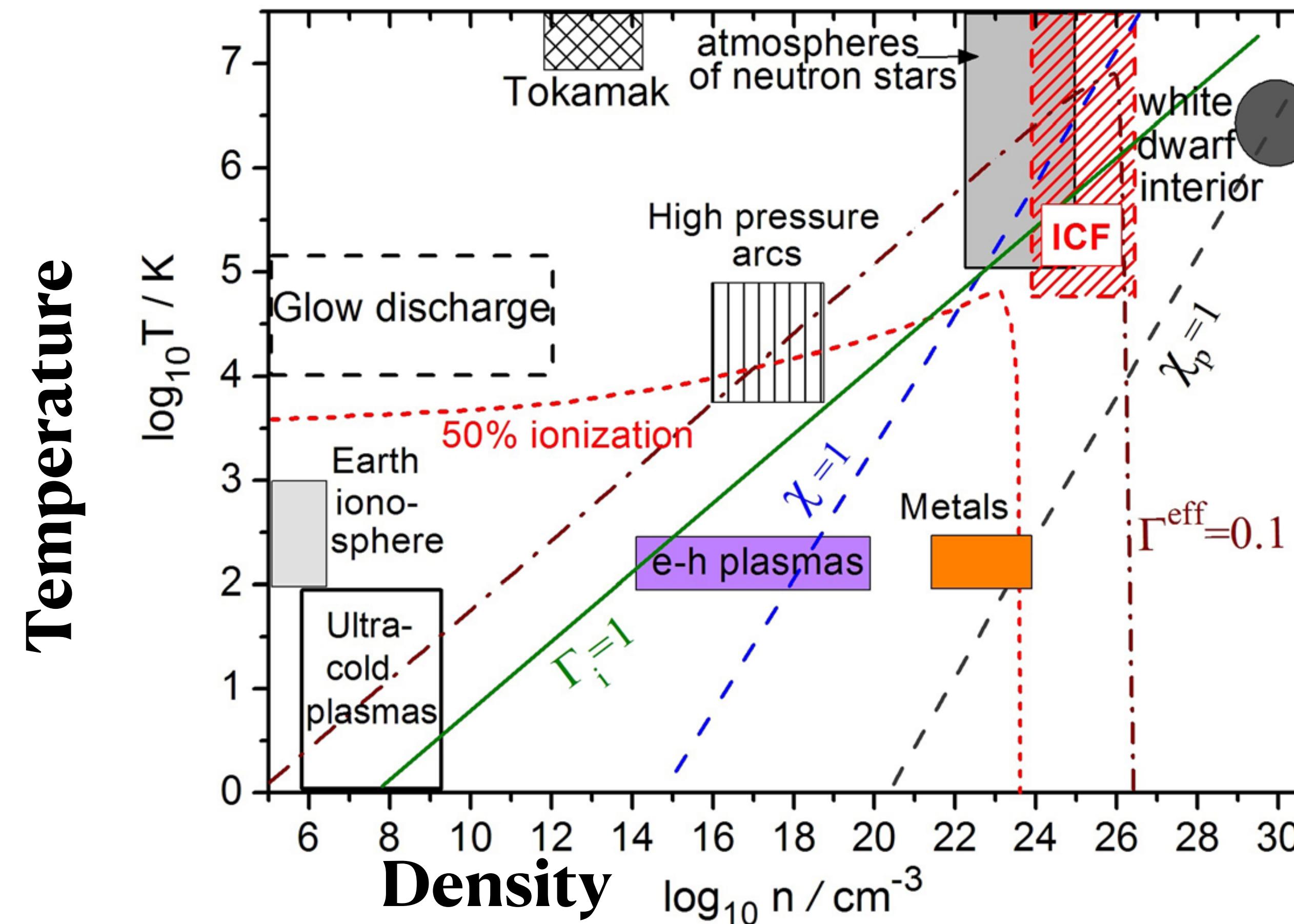
3

$$F = E - TS$$

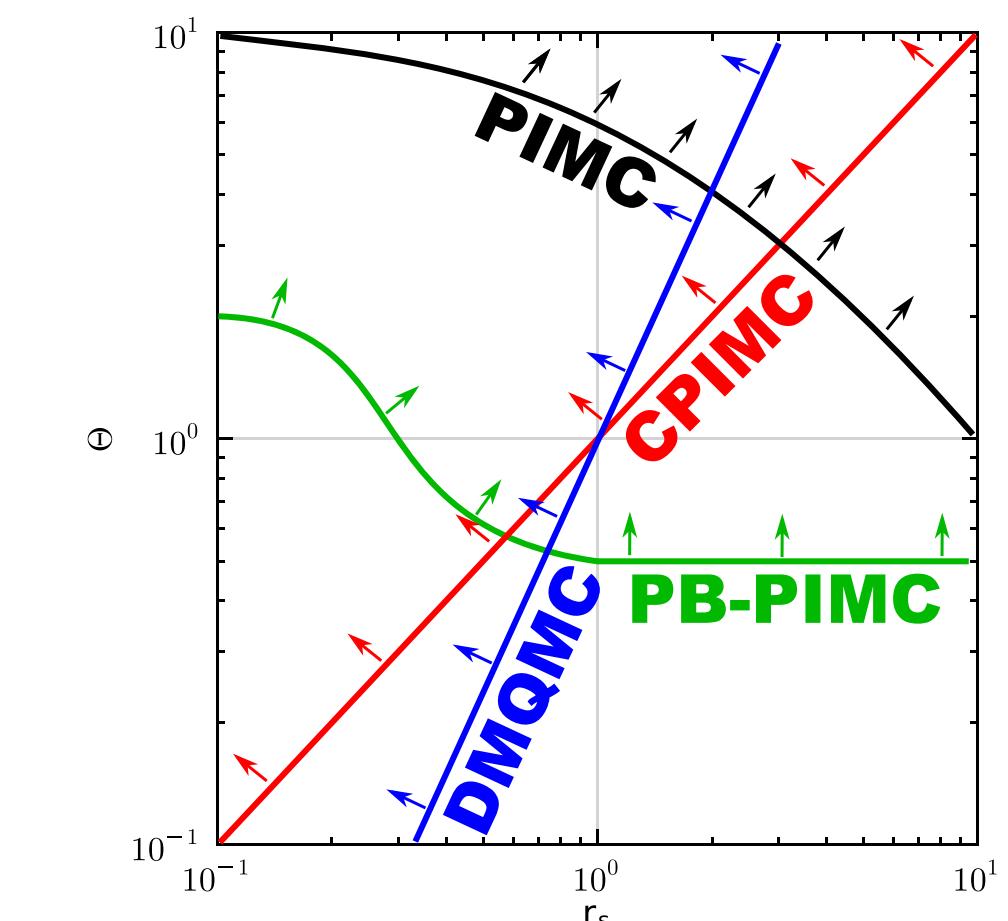
Nature's cost function

Ab-initio study of quantum matters at T>0

$$H = - \sum_i \frac{\hbar^2}{2m_e} \nabla_i^2 - \sum_I \frac{\hbar^2}{2m_I} \nabla_I^2 - \sum_{I,i} \frac{Z_I e^2}{|R_I - r_i|} + \frac{1}{2} \sum_{i \neq j} \frac{e^2}{|r_i - r_j|} + \frac{1}{2} \sum_{I \neq J} \frac{Z_I Z_J e^2}{|R_I - R_J|}$$



Quantum Monte Carlo methods
are limited by the “sign problem”



Dornheim et al, Phys. Plasmas '17
Bonitz et al, Phys. Plasmas '20

Classical world

Probability density p

Kullback-Leibler divergence

$$\mathbb{KL}(p || q)$$

Variational free-energy

$$F = \int dx \left[\frac{1}{\beta} p(x) \ln p(x) + p(x) E(x) \right]$$

Quantum world

Density matrix ρ

Quantum relative entropy

$$S(\rho || \sigma)$$

Variational free-energy

$$F = \frac{1}{\beta} \text{Tr}(\rho \ln \rho) + \text{Tr}(\rho H)$$

Variational density matrix with *two* generative models

$$\min F[\rho] = k_B T \operatorname{Tr}(\rho \ln \rho) + \operatorname{Tr}(H\rho)$$

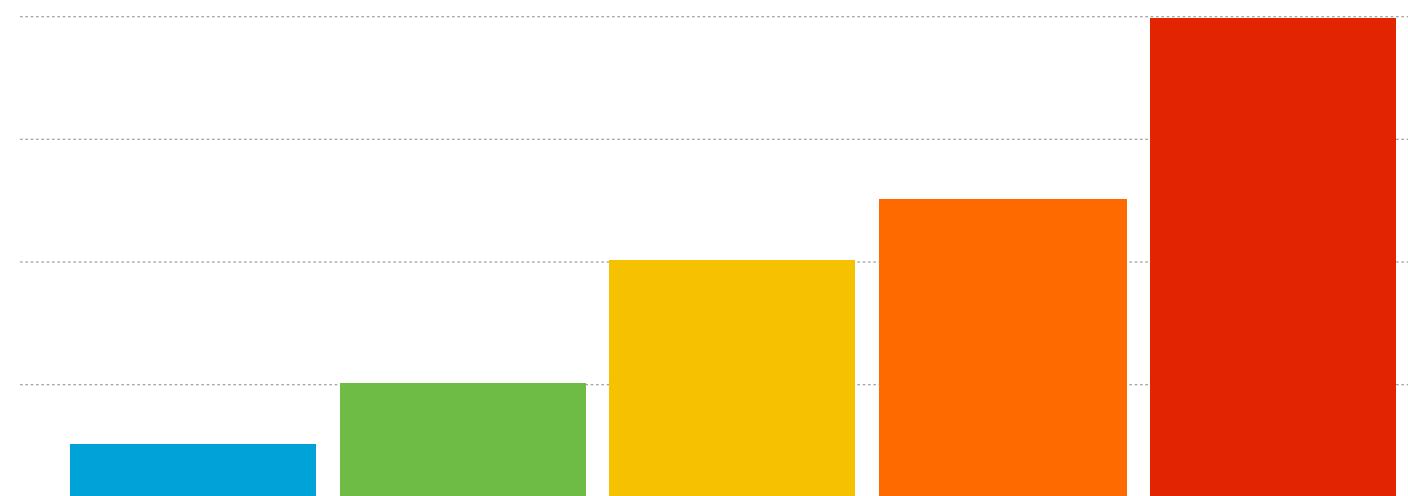
Gibbs–Bogolyubov–Feynman–
Delbrück–Molière

$$\text{s.t. } \operatorname{Tr}\rho = 1 \quad \rho > 0 \quad \rho^\dagger = \rho \quad \langle X | \rho | X' \rangle = (-)^{\mathcal{P}} \langle \mathcal{P}X | \rho | X' \rangle$$

$$\rho = \sum_n p_n |\Psi_n\rangle\langle\Psi_n|$$

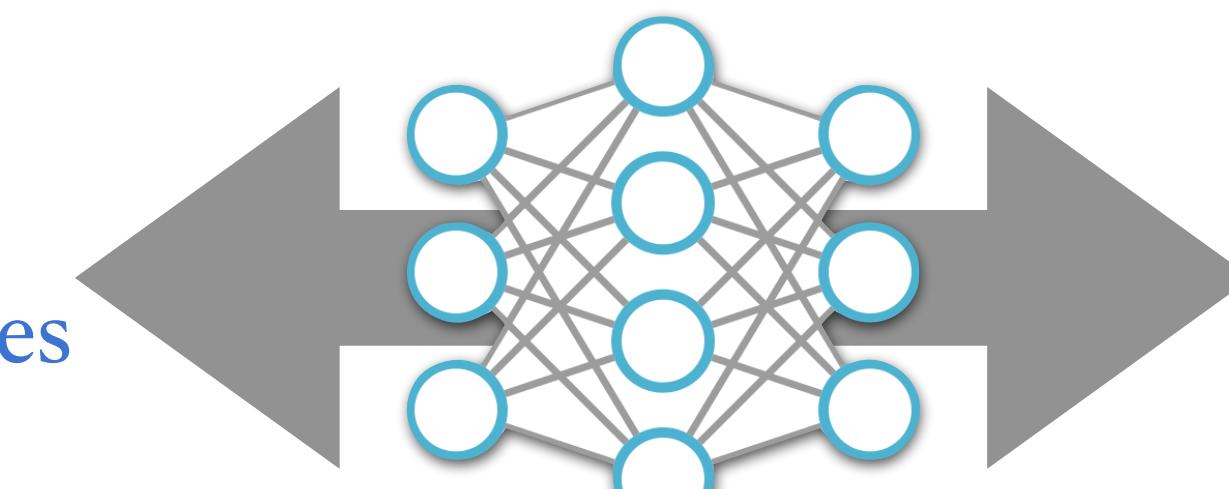
Xie, Zhang, LW, JML '22

Classical probability p_n



Discrete probabilistic models
e.g. an autoregressive model

particle
coordinates



$\sqrt{\text{Normalizing flow}}$

c.f. Cranmer et al, 1904.05903, Saleh et al, 2308.16468

Point Transformations in Quantum Mechanics

BRYCE SELIGMAN DEWITT*

Ecole d'Eté de Physique Théorique de l'Université de Grenoble, Les Houches, Haute Savoie, France

(Received September 14, 1951)

An isomorphism is shown to exist between the group of point transformations in classical mechanics and a certain subgroup of the group of all unitary transformations in quantum mechanics. This isomorphism is

The unitary representations of the point-transformation group may be obtained by determining the infinitesimal generators of the group. An infinitesimal point transformation may be expressed in the form

$$x'^i = x^i + \epsilon \Lambda^i(x), \quad (3.7)$$

$$p'_i = p_i - \frac{1}{2} \epsilon [(\partial/\partial x^i) \Lambda^j(x), p_j]_+, \quad (3.8)$$

Coordinate transformation induces a unitary $e^{\frac{i}{2}[\Lambda^i(x), p_i]_+}$

Point Transformations and the Many Body Problem*

M. EGER† AND E. P. GROSS

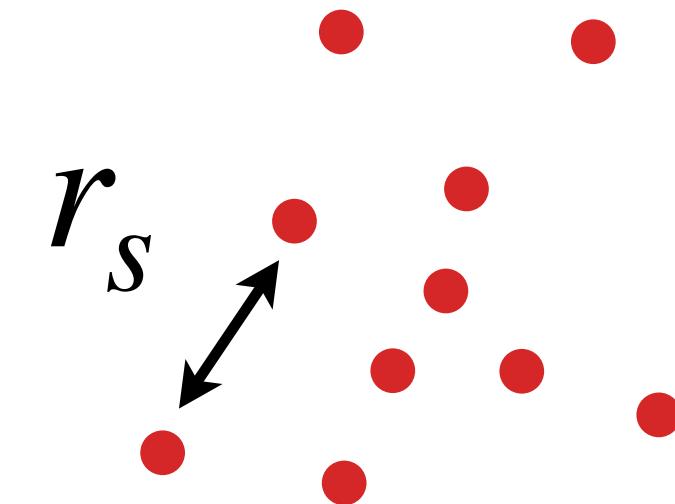
Brandeis University, Waltham, Massachusetts

An investigation is made of possible uses of many dimensional coordinate transformations in the quantum many-body problem. The transformed Hamiltonian is quadratic in the momenta with a space dependent metric. The original potential energy undergoes alteration and an additional "metric" potential energy appears. A relatively complete analysis of the transformed original potential is made, and the coordinate transformation can be used to suppress undesirable features of the original potential. For bosons one can attempt to directly map a complete set of noninteracting states onto approximate eigenstates of the system with interactions. Contact is made with a theory of weakly interacting bosons. In the general case it emerges that a given transformation uniquely fixes all the spatial correlation functions, which can be explicitly computed. The extended point transform can then be used as a link between diverse experimental quantities. The full use of the transformation to compute from first principles requires adequate approximations to the Jacobian and the inverse transform. These problems are not studied.

✓Normalizing flow
materialize this dream

Example: uniform electron gas

$$H = - \sum_{i=1}^N \frac{\hbar^2 \nabla_i^2}{2m} + \sum_{i < j} \frac{e^2}{|x_i - x_j|}$$

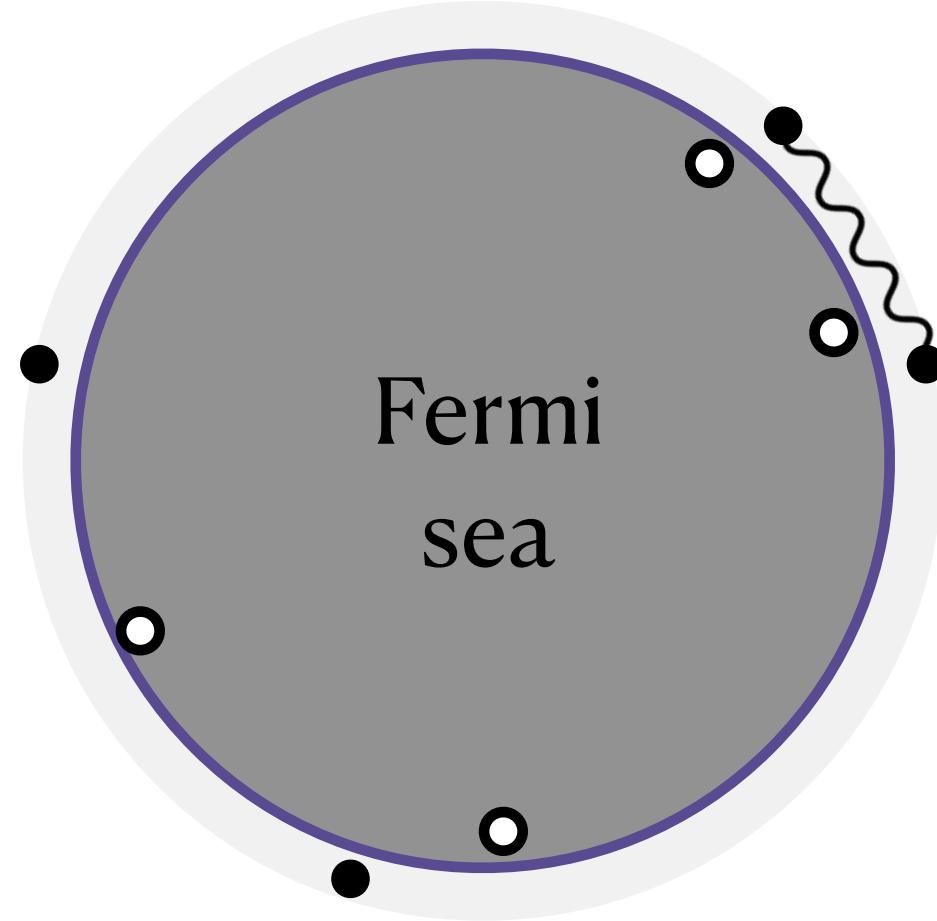


Fundamental model for metals ($2 < r_s < 6$)
Fermi liquid despite of non-perturbative r_s

$$E_c[n] = \int d^3r n(\epsilon_c^{\text{ueg}} + \dots)$$

Input to the density
functional theory calculations

Deep generative models for the variational density matrix



Low-energy excited
states are labeled in
the same way as the
ideal Fermi gas

$$K = \{k_1, k_2, \dots, k_N\}$$

$$\rho = \sum_K p(K) |\Psi_K\rangle\langle\Psi_K|$$

Normalized probability
distribution

Orthonormal
many-electron basis

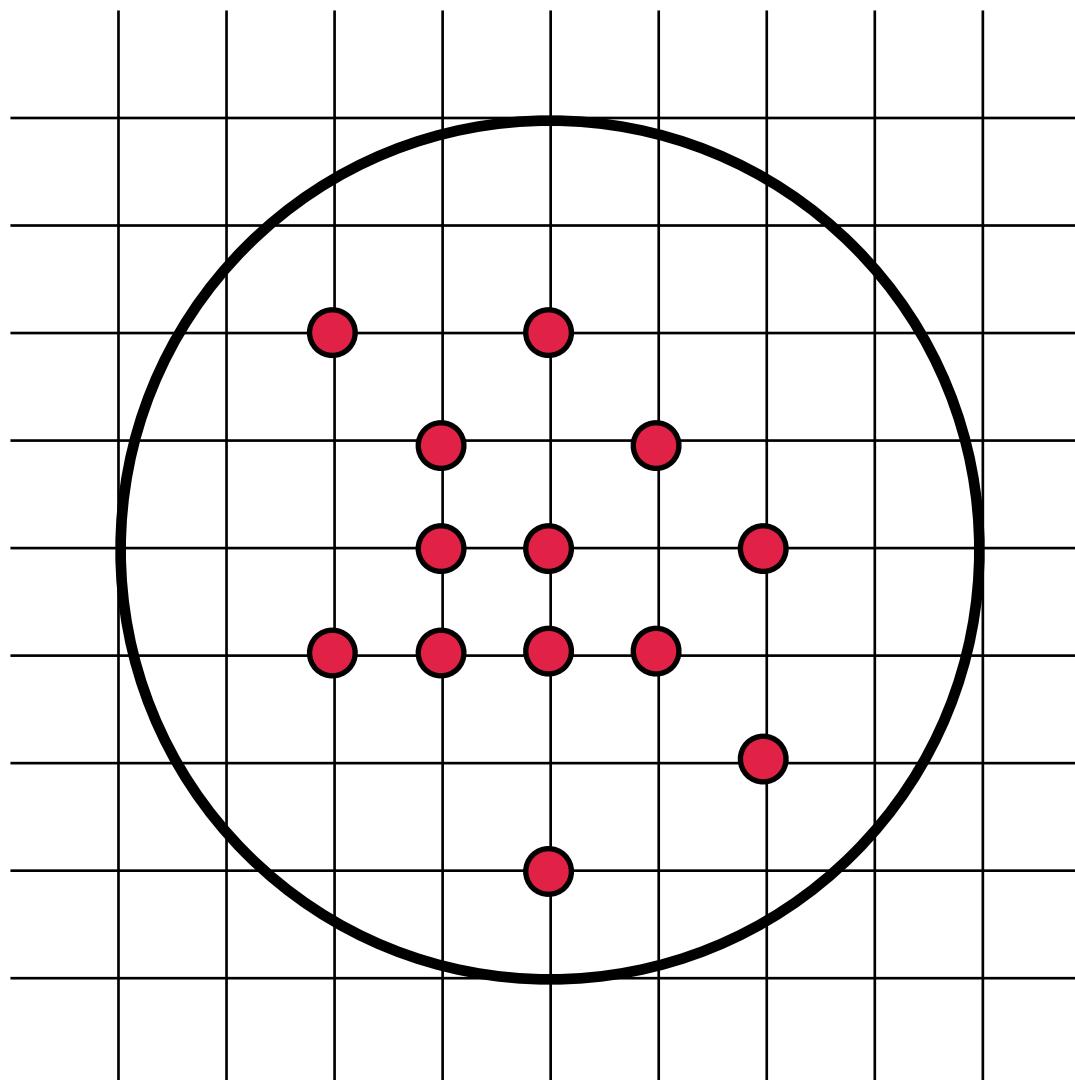
① $\sum_K p(K) = 1$

② $\langle\Psi_K | \Psi_{K'}\rangle = \delta_{K,K'}$

Imposing physics constraints into deep generative models

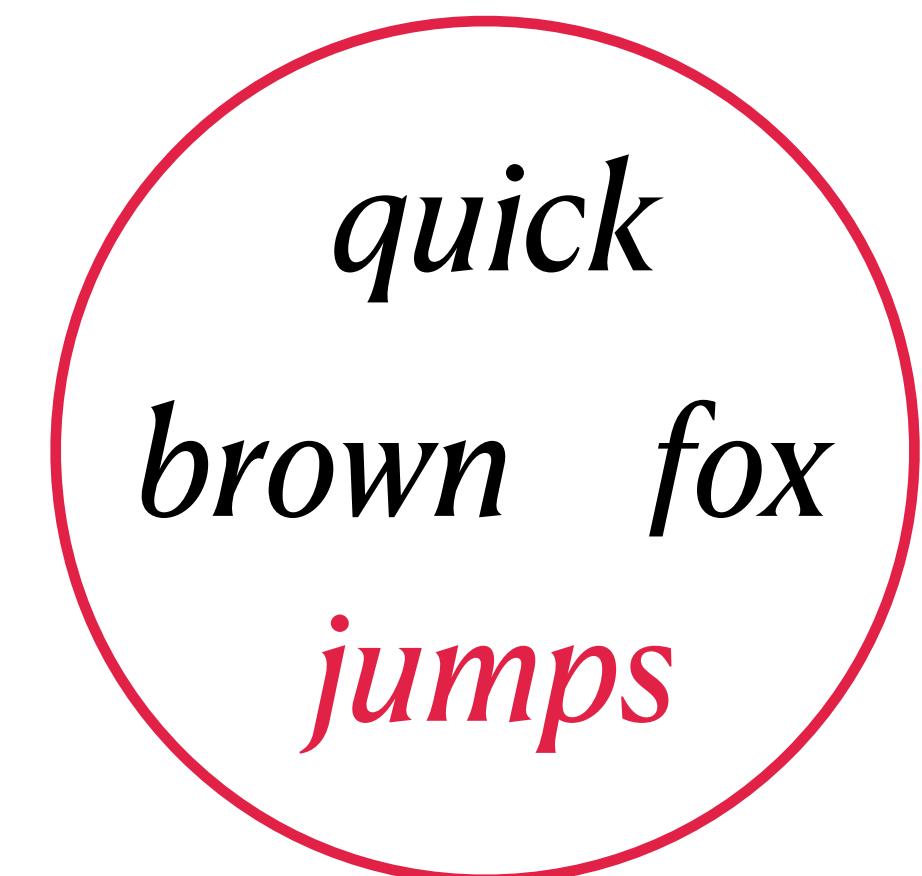
① Autoregressive model for $p(\mathbf{K})$

Fermionic
occupation
in k-space



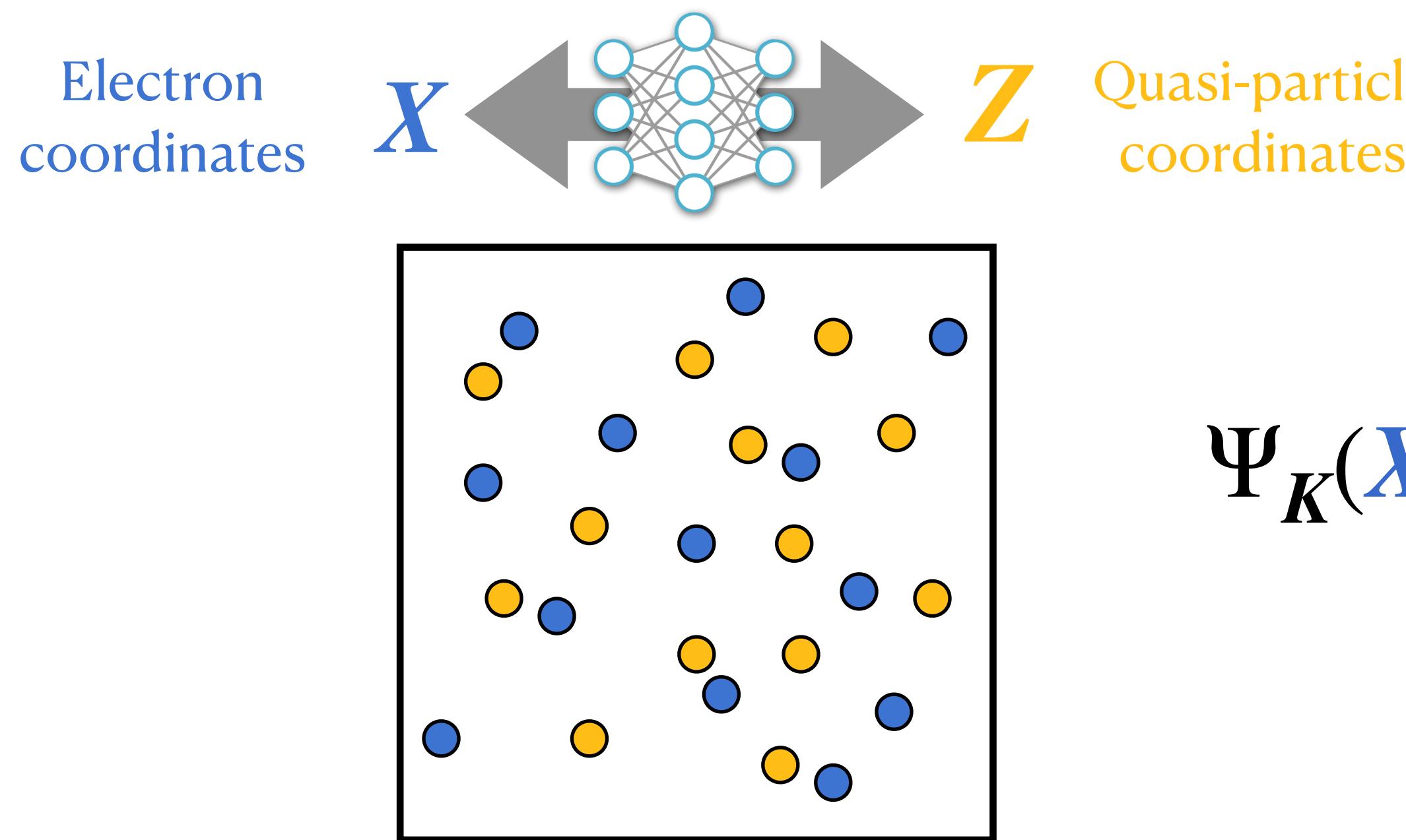
$$p(\mathbf{K}) = p(k_1)p(k_2 | k_1)p(k_3 | k_1, k_2)\dots$$

	Momentum distribution	Language
N	# of fermions	# of words
M	Momentum cutoff	Vocabulary
Space	$\binom{M}{N}$	M^N



Pauli exclusion: we are modeling a *set of words* with no repetitions and no order

$\sqrt{\text{Normalizing flow}}$ for $|\Psi_n\rangle$



$$\Psi_K(\mathbf{X}) = \frac{\det(e^{i\mathbf{k}_i \cdot \mathbf{z}_j})}{\sqrt{N!}} \cdot \left| \det \left(\frac{\partial \mathbf{Z}}{\partial \mathbf{X}} \right) \right|^{\frac{1}{2}}$$

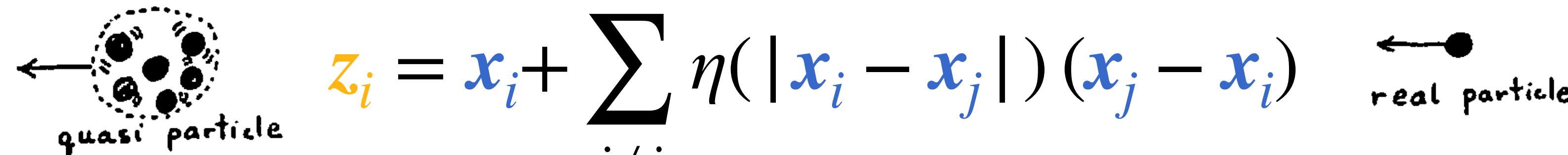
Orthonormal many-body states

Jacobian of the transformation

Fermion statistics: the flow should be permutation equivariant

we use FermiNet layer Pfau et al, 1909.02487, PRR '20

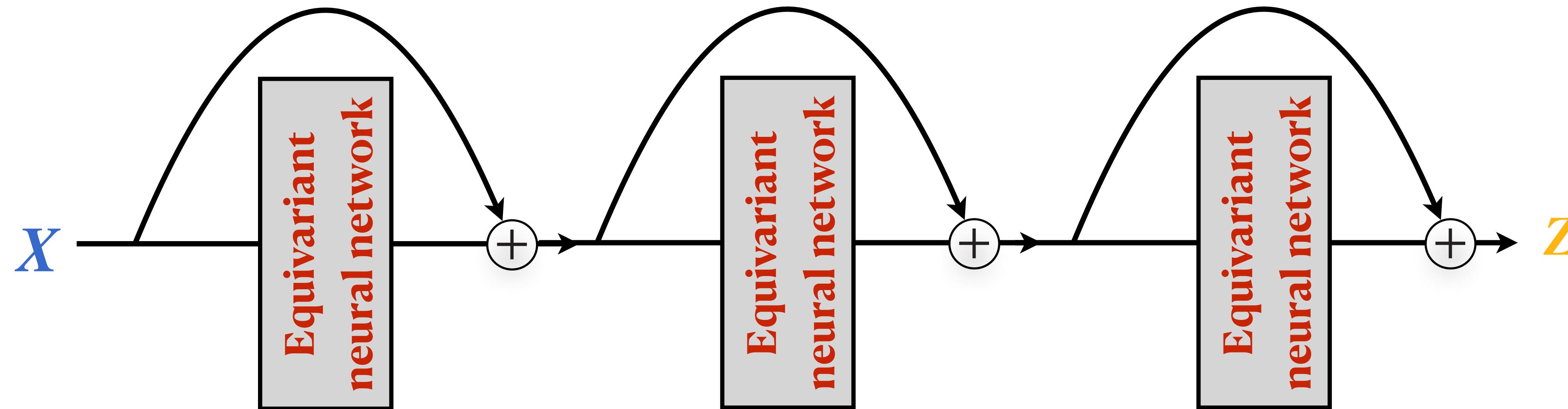
Feynman's backflow in the deep learning era



A diagram showing a cluster of black dots labeled "quasi particle" with an arrow pointing towards it. To its right is a single black dot labeled "real particle" with an arrow pointing away from it.

$$z_i = \mathbf{x}_i + \sum_{j \neq i} \eta(|\mathbf{x}_i - \mathbf{x}_j|) (\mathbf{x}_j - \mathbf{x}_i)$$

Feynman & Cohen 1956
wavefunction for liquid Helium



Iterative backflow → deep residual network → continuous normalizing flow



Fermi Flow

Xie, Zhang, LW, 2105.08644, JML '22

github.com/fermiflow

Continuous flow of electron density in a quantum dot



Fermi Flow

Xie, Zhang, LW, 2105.08644, JML '22

github.com/fermiflow

Continuous flow of electron density in a quantum dot

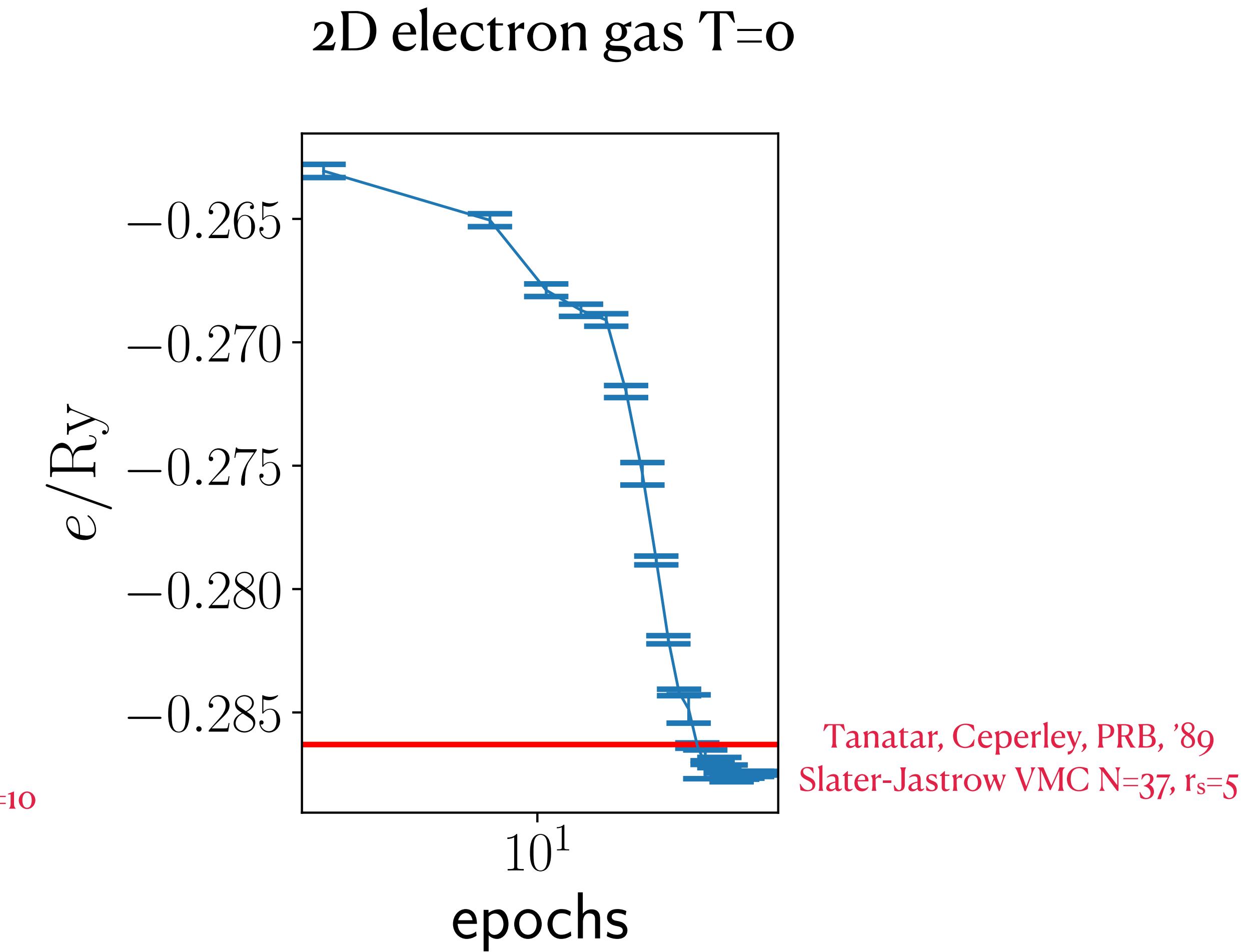
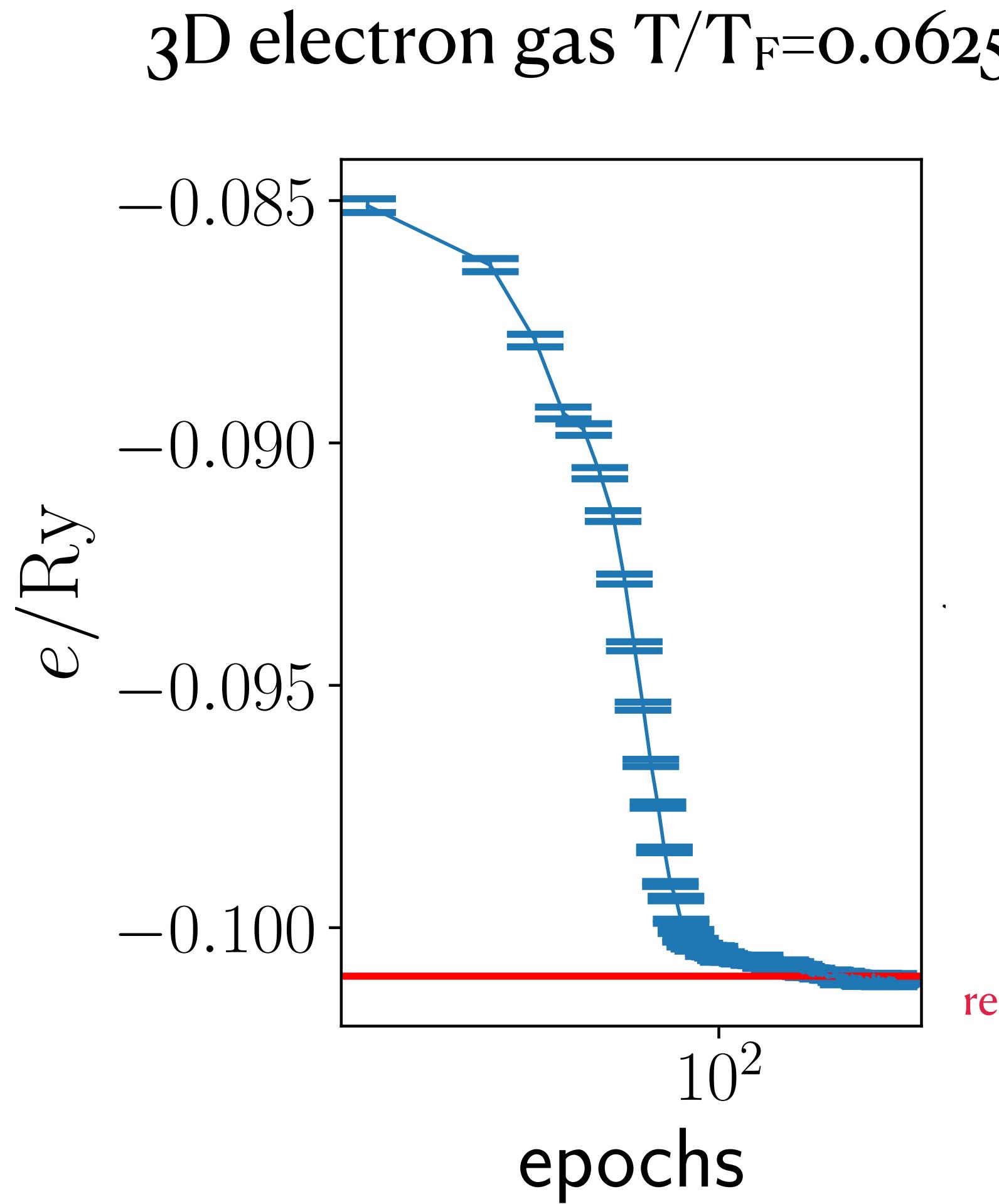
The objective function

$$F = \mathbb{E}_{K \sim p(K)} \left[k_B T \ln p(K) + \mathbb{E}_{X \sim |\langle X | \Psi_K \rangle|^2} \left[\frac{\langle X | H | \Psi_K \rangle}{\langle X | \Psi_K \rangle} \right] \right]$$

↓ ↓
Boltzmann Born
distribution probability

Jointly optimize $|\Psi_K\rangle$ and $p(K)$ to minimize the variational free energy

Benchmarks on spin-polarized electron gases



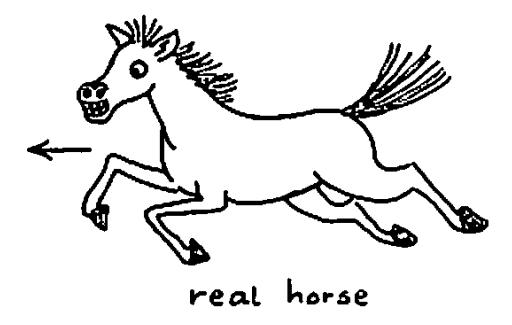
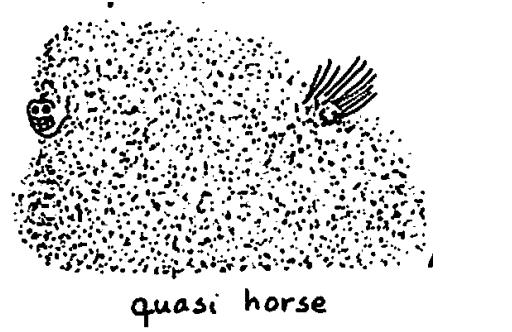
Application: m^* from low temperature entropy

Eich, Holzmann, Vignale, PRB '17

$$S = \frac{\pi^2 k_B}{3} \frac{m^*}{m} \frac{T}{T_F}$$

$$\Rightarrow \frac{m^*}{m} = \frac{s}{s_0}$$

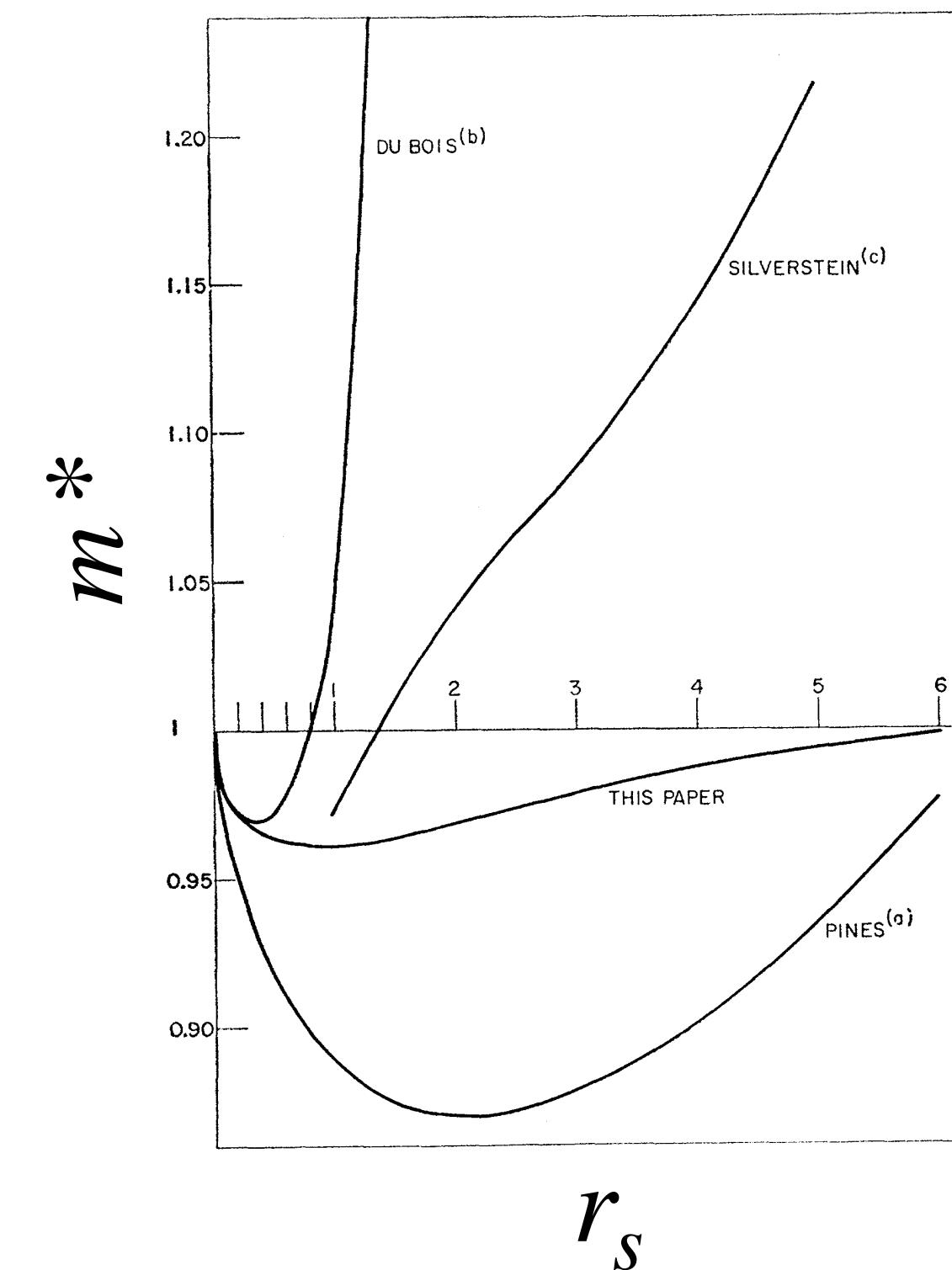
interacting electrons
noninteracting electrons



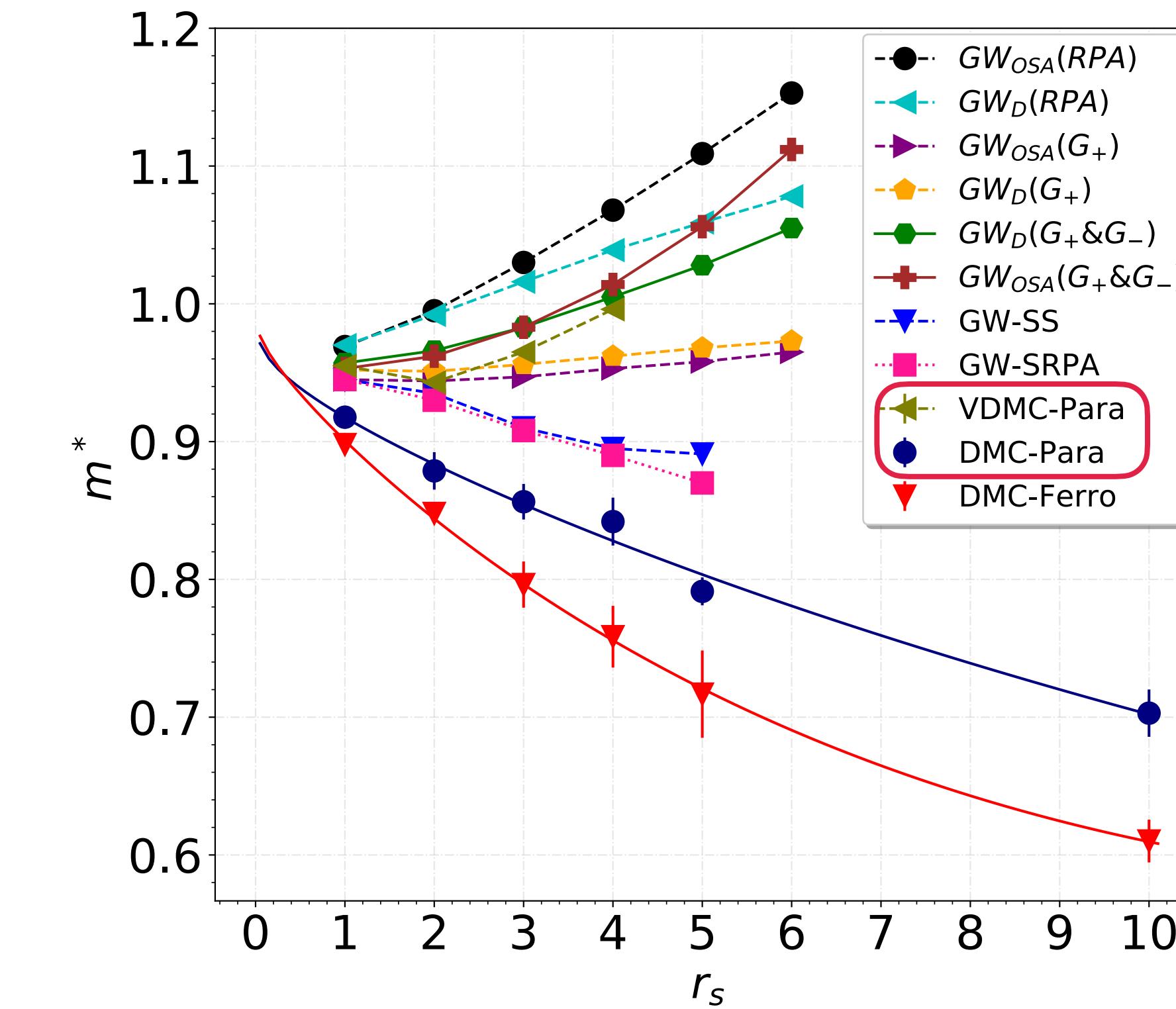
A fundamental quantity appears in nearly all physical properties of a Fermi liquid
There have been debates despite its fundamental role and long history of study

Quasi-particles effective mass of 3d electron gas

Hedin Phy. Rev. 1965



Azadi, Drummond, Foulkes, PRL 2021



>50 years of conflicting results !

Two-dimensional electron gas experiments

VOLUME 91, NUMBER 4

PHYSICAL REVIEW LETTERS

week ending
25 JULY 2003

Spin-Independent Origin of the Strongly Enhanced Effective Mass in a Dilute 2D Electron System

A. A. Shashkin,* Maryam Rahimi, S. Anissimova, and S.V. Kravchenko

Physics Department, Northeastern University, Boston, Massachusetts 02115, USA

V.T. Dolgopolov

Institute of Solid State Physics, Chernogolovka, Moscow District 142432, Russia

T. M. Klapwijk

Department of Applied Physics, Delft University of Technology, 2628 CJ Delft, The Netherlands

(Received 13 January 2003; published 24 July 2003)

$$m^*/m > 1$$



PRL 101, 026402 (2008)

PHYSICAL REVIEW LETTERS

week ending
11 JULY 2008

Effective Mass Suppression in Dilute, Spin-Polarized Two-Dimensional Electron Systems

Medini Padmanabhan, T. Gokmen, N. C. Bishop, and M. Shayegan

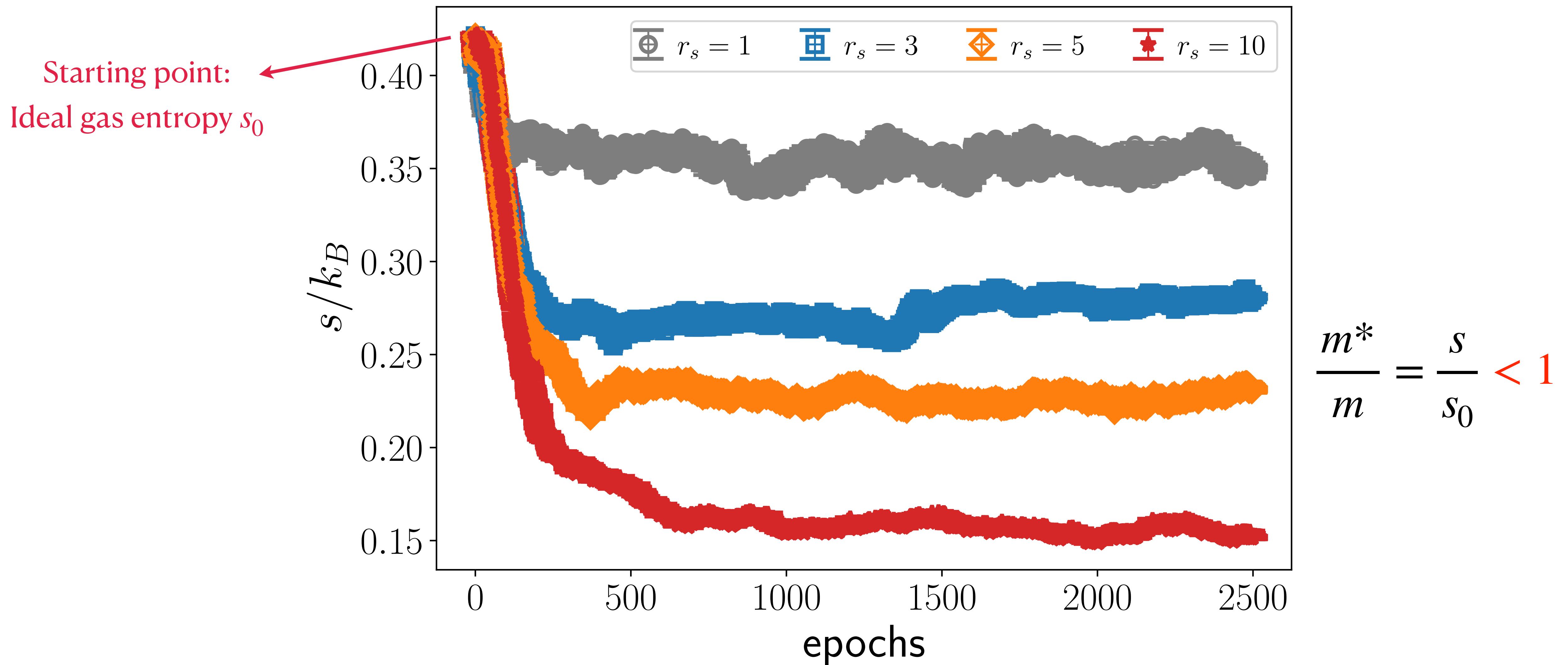
Department of Electrical Engineering, Princeton University, Princeton, New Jersey 08544, USA

(Received 19 September 2007; published 7 July 2008)

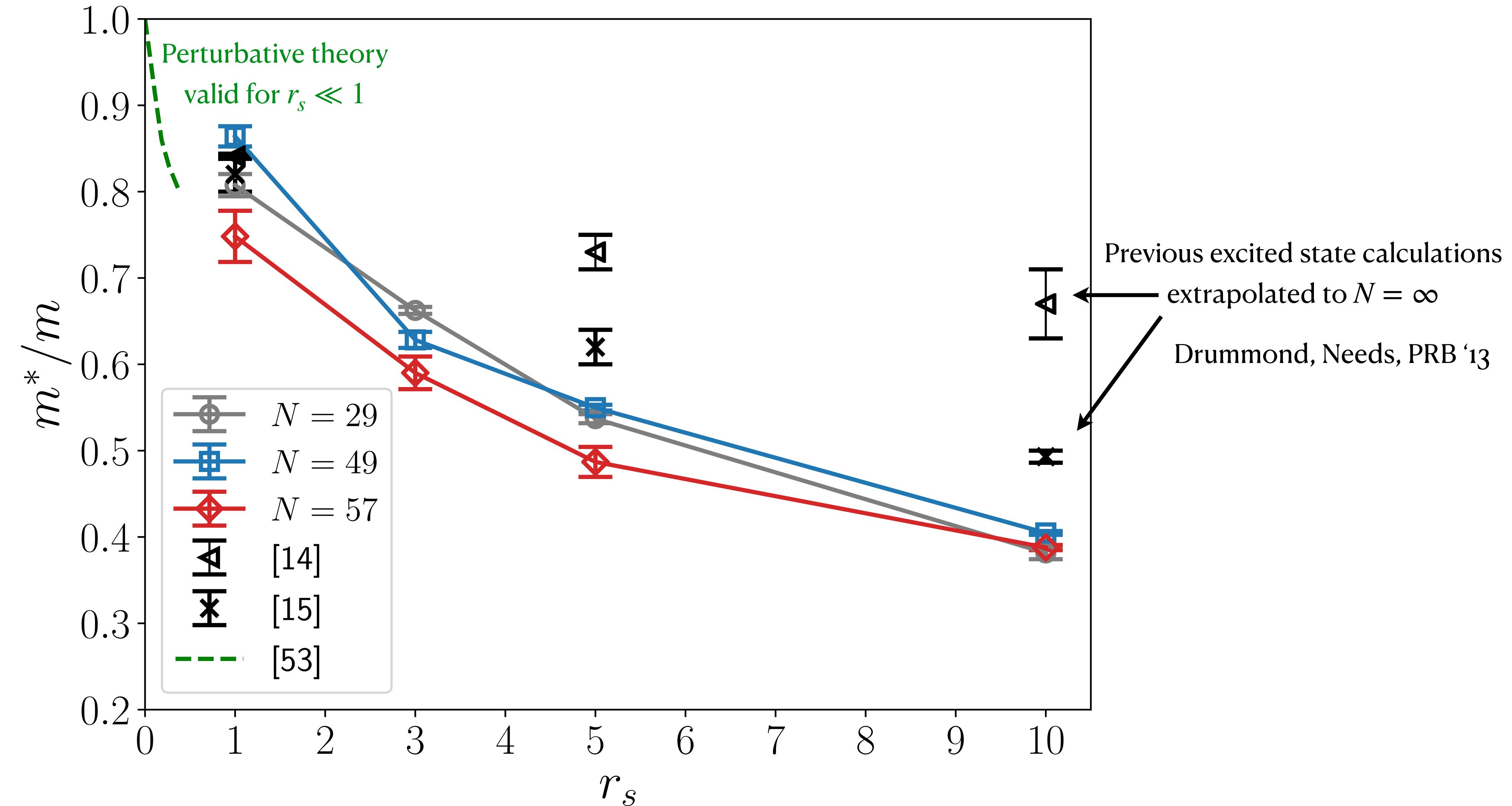
$$m^*/m < 1$$

Layer thickness, valley, disorder, spin-orbit coupling...

37 spin-polarized electrons in 2D @ T/T_F=0.15



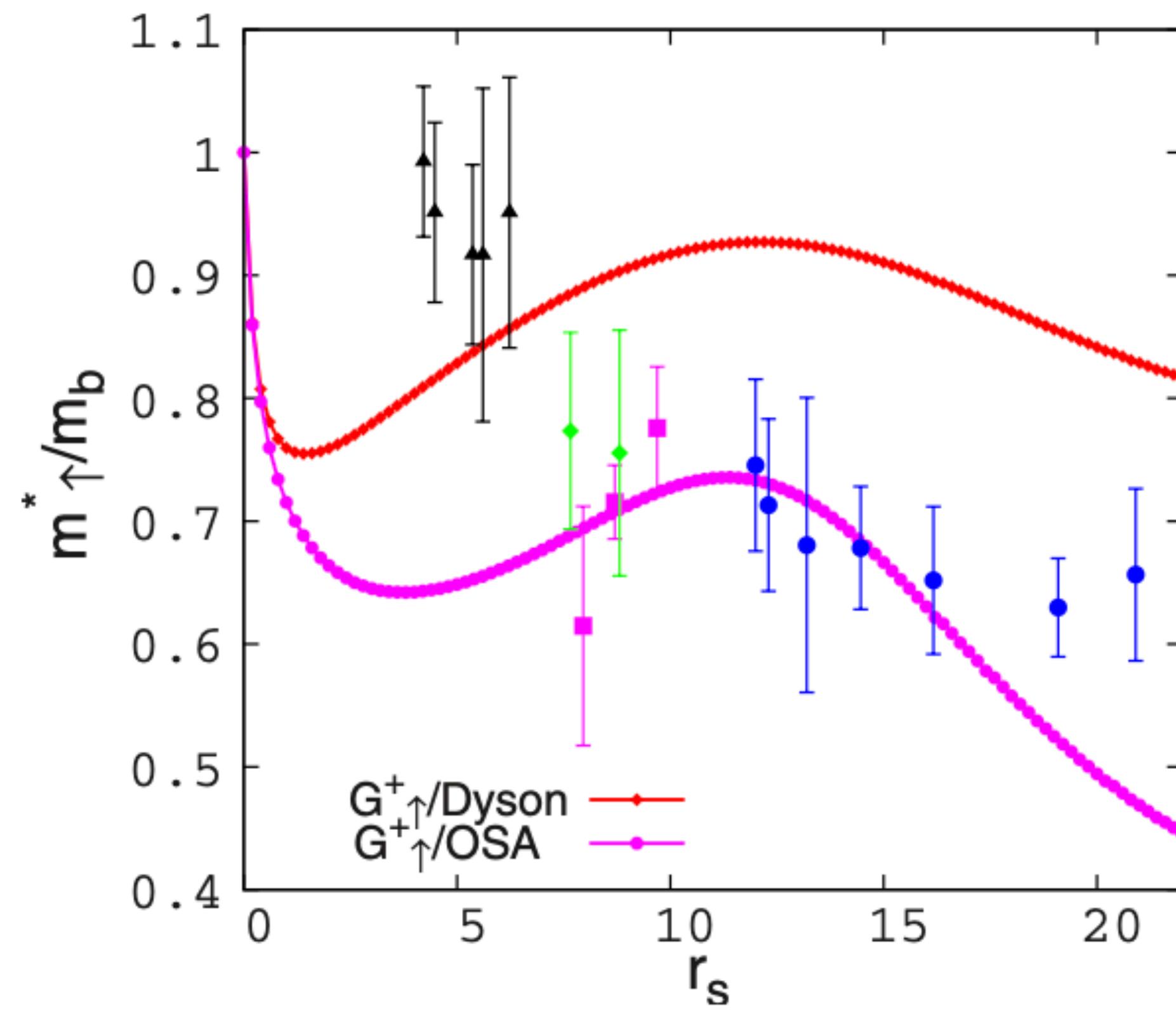
Effective mass of spin-polarized 2DEG



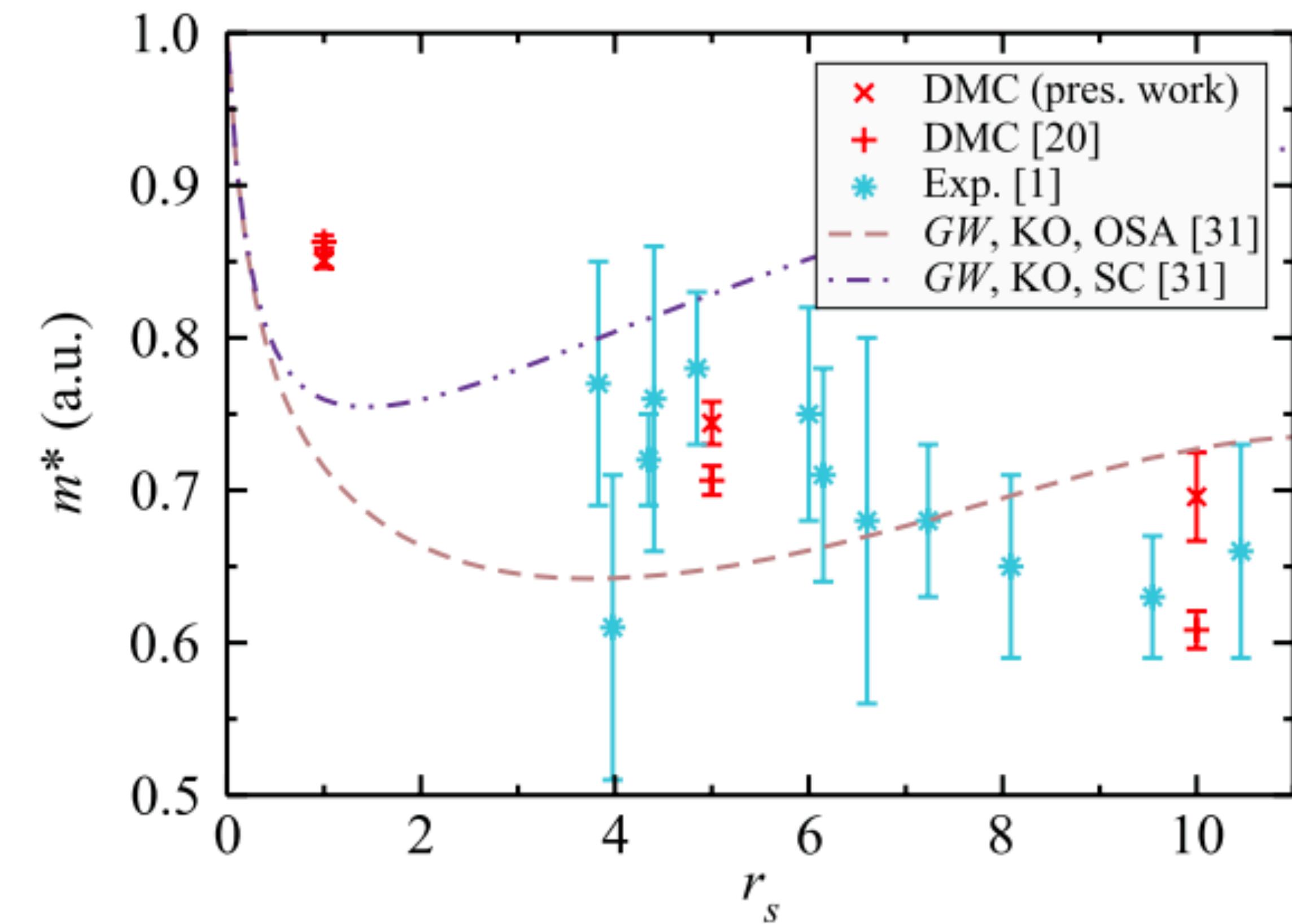
More pronounced suppression of m^* in the low-density strong-coupling region

Experiments on spin-polarized 2DEG

Asgari et al, PRB '09



Drommond, Needs, PRB'13



Quantum oscillation experiments
Padmanabhan et al, PRL '08
Gokmen et al, PRB '09

Entropy measurement of 2DEG

ARTICLE

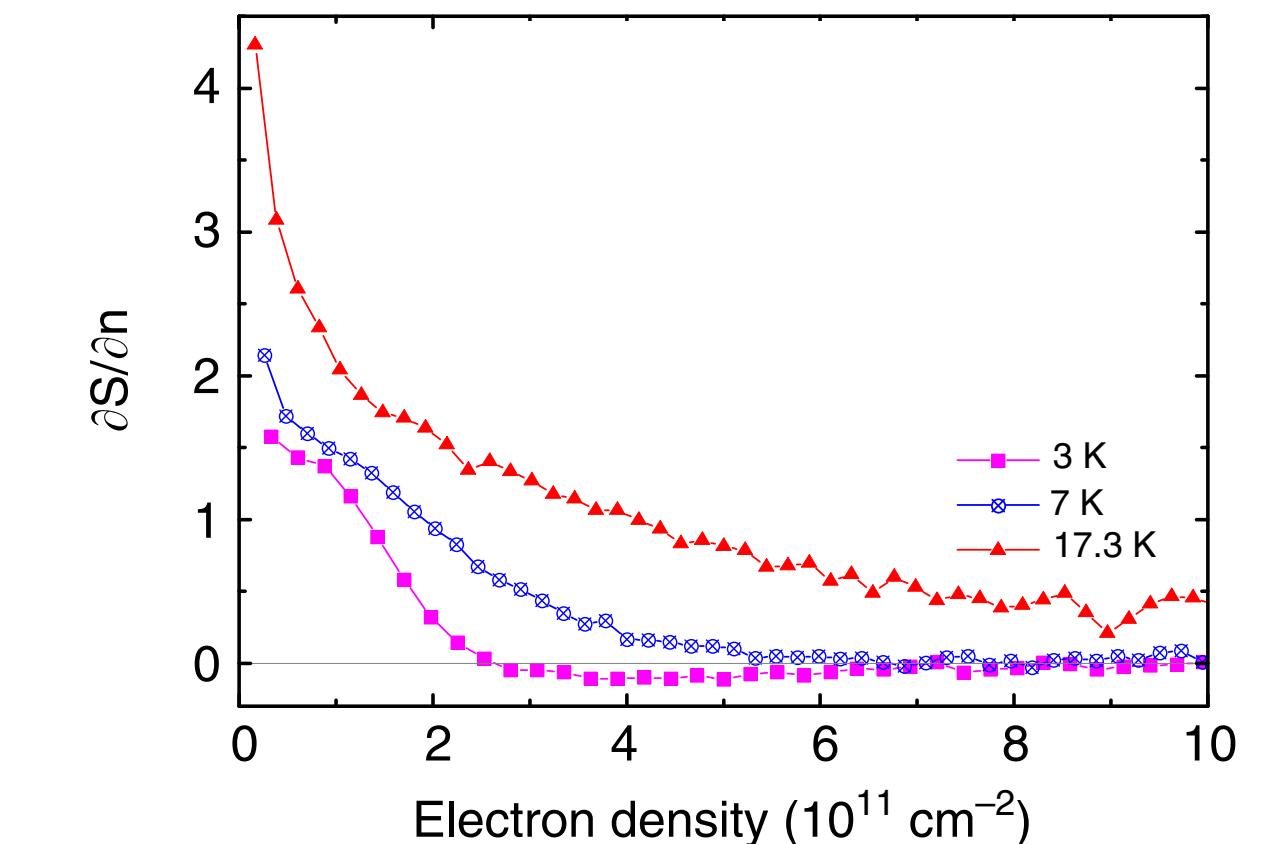
Received 16 May 2014 | Accepted 27 Apr 2015 | Published 23 Jun 2015

DOI: [10.1038/ncomms8298](https://doi.org/10.1038/ncomms8298)

Strongly correlated two-dimensional plasma explored from entropy measurements

A.Y. Kuntsevich^{1,2}, Y.V. Tupikov³, V.M. Pudalov^{1,2} & I.S. Burmistrov^{2,4}

Maxwell relation $\left(\frac{\partial S}{\partial n}\right)_T = - \left(\frac{\partial \mu}{\partial T}\right)_n$



Next, directly compare computed entropy with the experiment

FAQs

Where to get training data ?

No training data. Data are self-generated from the generative model.

How do we know it is correct ?

Variational principle: lower free-energy is better.

Do I understand the “black box” model ?

- a) I don't care (as long as it is sufficiently accurate).
- b) $\ln p(K)$ contains the Landau energy functional

$Z \leftrightarrow X$ illustrates adiabatic continuity.

$$E[\delta n_k] = E_0 + \sum_k \epsilon_k \delta n_k + \frac{1}{2} \sum_{k,k'} f_{k,k'} \delta n_k \delta n_{k'}$$

“Using AI to accelerate scientific discovery” Demis Hassabis, co-founder and CEO of DeepMind 2021

What makes for a suitable problem?

1

Massive combinatorial
search space

2

Clear objective function
(metric) to optimise
against

3

Either lots of data
and/or an accurate and
efficient simulator

Why now ?

Variational free-energy is a **fundamental principle** for $T > 0$ quantum systems

However, it was under-exploited for solving practical problems
(mostly due to intractable entropy for nontrivial density matrices)

Now, it has became possible by integrating recent advances in
generative models

The Universe as a generative model

$$S = \int d^4x \sqrt{-g} \left[\frac{m_p^2}{2} R - \frac{1}{4} F_{\mu\nu}^a F_a^{\mu\nu} + i \bar{\psi}^i D_\mu \psi^i + \left(\bar{\psi}_L^i V_{ij} \not{D} \psi_R^j + h.c. \right) - |\not{D}_\mu \Phi|^2 - V(\Phi) \right]$$



Thank you!

Discovering physical laws: **learning** the action
Solving physical problems: **optimizing** the action

2.23

Overview

3.2

Machine learning practices

3.9

A hitchhiker's guide to deep learning

3.16

Research projects hands-on

3.23

Symmetries in machine learning

3.30

Differentiable programming

4.6

Generative models-I

4.13

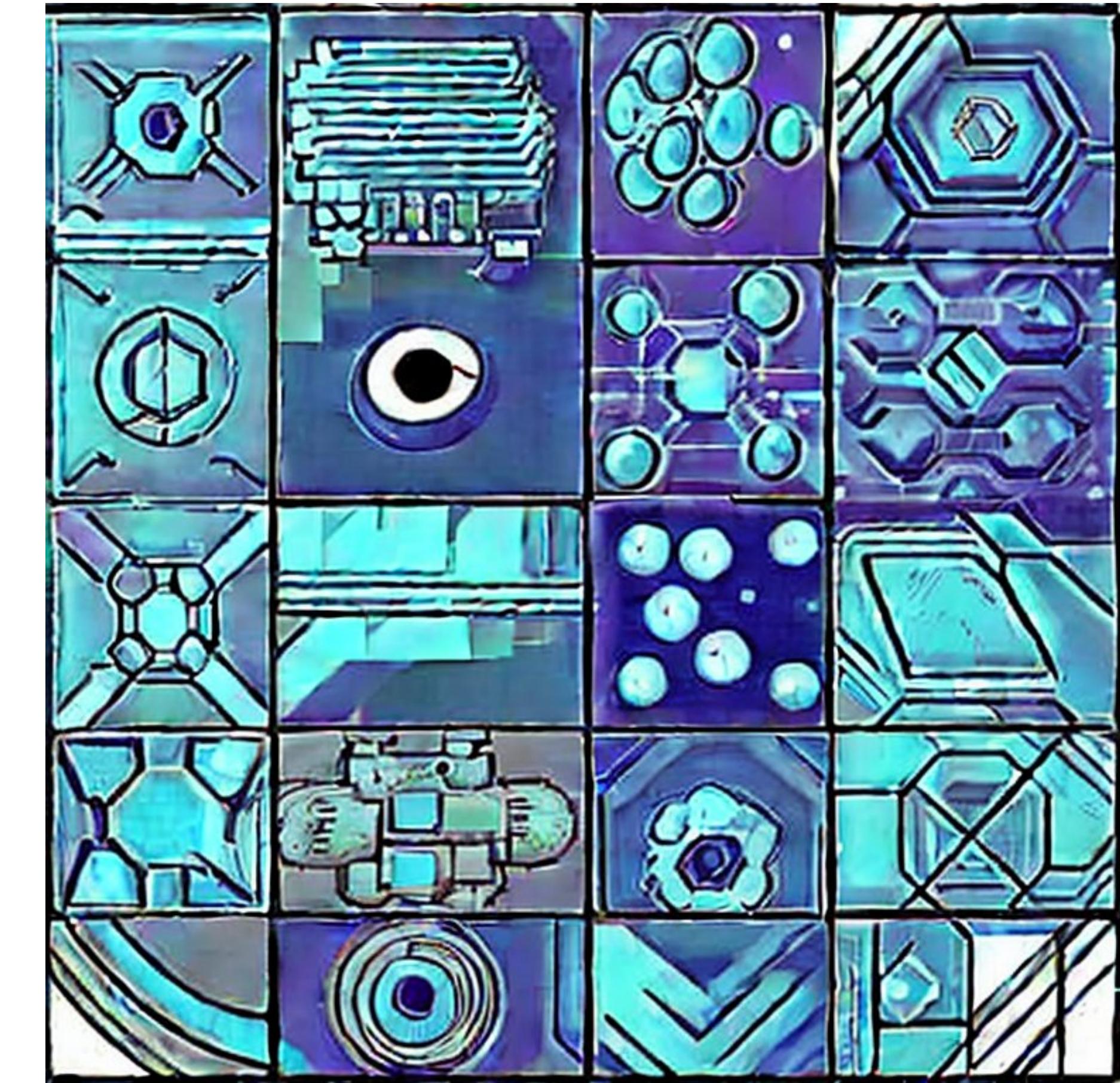
Generative models-II

4.20

Research projects presentation

4.27

AI for science: why now ?



Machine learning for physicists

<https://github.com/wangleiphy/ml4p>

Machine Learning: Science and Technology

Focus on Generative AI in Science

Guest Editors

Juan Felipe Carrasquilla, *Vector Institute, Canada*

Stephen R. Green, *University of Nottingham, UK*

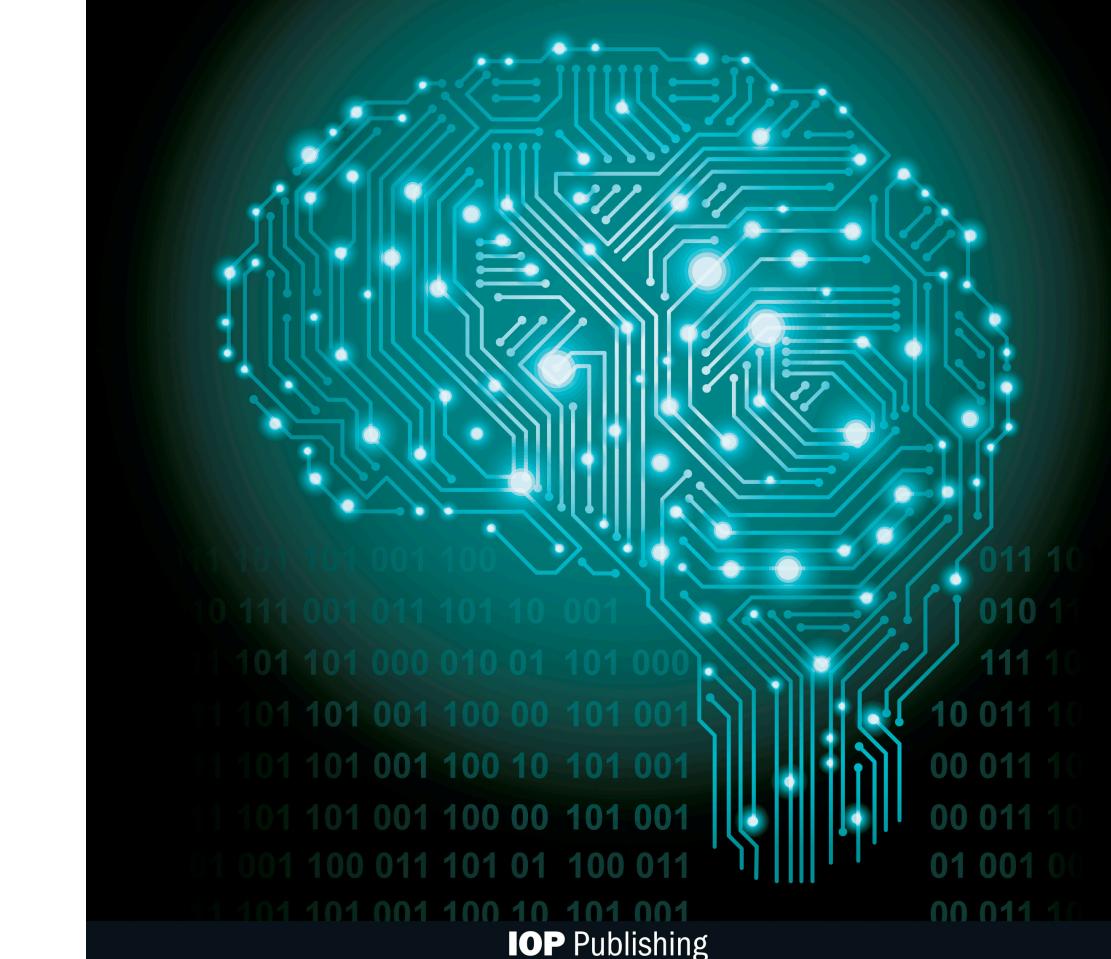
Lei Wang, *Institute of Physics, CAS, China*

Linfeng Zhang, *DP Technology/AI for Science Institute, China*

Pan Zhang, *Institute of Theoretical Physics, CAS, China*

MACHINE
LEARNING
Science and Technology

iopscience.org/mlst



<https://iopscience.iop.org/collections/mlst-230424-207>